

Informed baseline subtraction of proteomic mass spectrometry data aided by a novel sliding window algorithm

Tyman E. Stanford^{a,†}, Christopher J. Bagley , and Patty J. Solomon^a

^aSchool of Mathematical Sciences, The University of Adelaide

[†]To whom correspondence should be addressed: School of Mathematical Sciences, The University of Adelaide, Ingkarni Wardli Building, North Terrace, Adelaide 5005, Australia, Email:

tyman.stanford@adelaide.edu.au

Abstract

Background Proteomic matrix-assisted laser desorption/ionisation (MALDI) linear time-of-flight (TOF) mass spectrometry (MS) may be used to produce protein profiles from biological samples with the aim of discovering biomarkers for disease or discrimination of disease states. The raw protein profiles suffer from several sources of bias or systematic variation, known as batch effects, which need to be removed before meaningful downstream analysis of the data can be undertaken. An early pre-processing step is baseline subtraction, which is the removal of non-peptide signal from the spectra. Baseline subtraction is complicated by each spectrum having, on average, wider peaks for peptides with higher mass-to-charge ratios (m/z). Additionally, the trial-and-error process of optimising the baseline subtraction input arguments is time-consuming and error-prone. We present an analytical pipeline to overcome these current difficulties.

Methods Current best practice baseline subtraction is performed by partitioning the spectra into smaller regions. The baseline subtraction method is then applied with constant and optimised input arguments within each region. We propose a new approach which transforms the m/z -axis to remove the relationship between peptide mass and peak width. Our preferred baseline subtraction method of the top-hat operator employs fast sliding window algorithms such as the line segment algorithm which cannot be applied to unevenly spaced data. We have also developed a novel ‘continuous’ line segment algorithm to efficiently operate on unevenly spaced data. To reduce the need for user input and the possibility of user error, we additionally present an input-free algorithm to estimate peak widths on the transformed m/z scale and thus the required sliding window widths for the top-hat operator. The methods are validated using six publicly available proteomic MS datasets.

Results The automated baseline subtraction method was deployed on each dataset using six different m/z -axis transformations. The resulting baseline subtracted signal was compared to the gold-standard piecewise baseline subtracted signal. Optimality of the m/z -axis transformation when using the automated baseline subtraction pipeline was assessed quantitatively using the mean absolute scaled error (MASE). Several of the transformations investigated were able to reduce, if not entirely remove, the peak width and peak location relationship. The best performing transformations achieved automated baseline subtractions very similar to the gold-standard. The proposed novel ‘continuous’ line segment algorithm is shown to far outperform naive sliding window algorithms with regard to the computational time required, on both real and simulated unevenly spaced MALDI TOF-MS data. The improvement observed in the time required to compute baseline subtraction on the six MALDI TOF-MS datasets was at least four-fold and at least an order of magnitude on many simulated datasets.

Conclusions The new pipeline presented here for performing baseline subtraction has a number of advantages over currently available methods. These advantages are: informed and data specific input arguments for baseline subtraction methods, the avoidance of time-intensive and subjective piecewise baseline subtraction, and the ability to automate baseline subtraction completely. Moreover, individual steps can be adopted as stand-alone routines. For example, the algorithm to automatically estimate peak widths can be used to dynamically calculate initial baseline subtraction method input arguments for subsequent user refinement for any given dataset. The proposed automated pipeline produced near-optimal baseline subtraction when compared to the current gold-standard of piecewise baseline subtraction method.

Keywords: mathematical morphology, top-hat operator, line segment algorithm, mass spectrometry, baseline subtraction, pre-processing, matrix-assisted laser desorption/ionization, time-of-flight, unevenly spaced data

Background

Discovery of protein biomarkers by mass spectrometry

Protein biomarkers are proteins or protein fragments that serve as markers of a disease or condition biomarkers [1, 2] by virtue of their altered relative abundance in the disease state versus the healthy condition. Matrix-assisted laser desorption/ionisation (MALDI) linear time-of-flight (TOF) mass spectrometry (MS) is a widely used technology for biomarker discovery as it can create a representative profile of polypeptide expression from biological samples. These profiles are displayed as points of polypeptide abundance (intensity; the y -axis) for a range of mass-to-charge values (m/z ; the x -axis). Each spectrum is an array of positive intensity values for discretely measured m/z values, but the profile is typically displayed on a continuous

scale. MALDI TOF-MS spectra are typically limited to polypeptides less than 30 kilo Daltons although there is no theoretical upper limit [3]. Numerous biomarkers using MALDI TOF-MS have been identified to date [3, 4].

Statistical analysis of the proteomic profiles for biomarker discovery cannot be undertaken without prior removal of noise and systematic bias present in the raw spectra. This removal is conducted through a series of steps known as pre-processing. Pre-processing generally consists of five steps to remove false signal as set out in Figure 1: signal smoothing, baseline subtraction, normalisation, peak detection and peak alignment. Signal smoothing and baseline subtraction are adjustments made to each spectrum individually (i.e., intra-spectrum pre-processing), while normalisation and peak alignment (after peak detection) are adjustments made to make each spectrum within an experiment comparable (i.e., inter-spectrum pre-processing).

Signal smoothing is the first step in pre-processing the data and aims to remove instrument-derived noise in the data and stochastic variation in the spectrum signal. Baseline subtraction then follows, which is the removal of the estimated ‘bed’ on which the spectral profile sits, composed of non-biological signal, e.g. chemical noise from ionised matrix. Normalisation is the third step in pre-processing. This has the aim of making the observed signals proportionate over the experiment; to correct for instrument variability and sample-ionisation efficiency that will influence the number of peptide ions reaching the detector. Peak detection is the fourth step, which is the detection of peak signal as peptide mass and intensity pairs. Finally, in the fifth step, the peaks are subject to peak alignment which adjusts for small drifts in m/z location which result from the calibration required for the TOF-MS system. This ensures that peptides common across spectra are recognised and compared at the same m/z value. Once the data have been pre-processed, analysis to detect potential biomarkers can be performed.

There are numerous freely available MS pre-processing packages. For example, in the R statistical software environment, **MALDIquant**, **PROcess** and **XCMS** are available [5, 6, 7, 8, 9]. Although we have set out the usual sequence of five data pre-processing steps, an optimal approach to pre-processing is not yet established and there is scope to improve current pre-processing methods and the order in which they are applied, to allow more reliable biomarker identification [10]. The present paper focuses on

Abbreviations used: AMASE: average mean absolute scaled error; CLSA: continuous line segment algorithm; Da: Dalton; EPCP: estimated peak coverage proportion; LOESS: Locally weighted scatterplot smoothing; LSA: line segment algorithm; MALDI: matrix-assisted laser desorption/ionization; MASE: mean absolute scaled error; MS: mass spectrometry; MSE: mean squared error; SNIP: sensitive nonlinear iterative peak; TOF: time of flight.

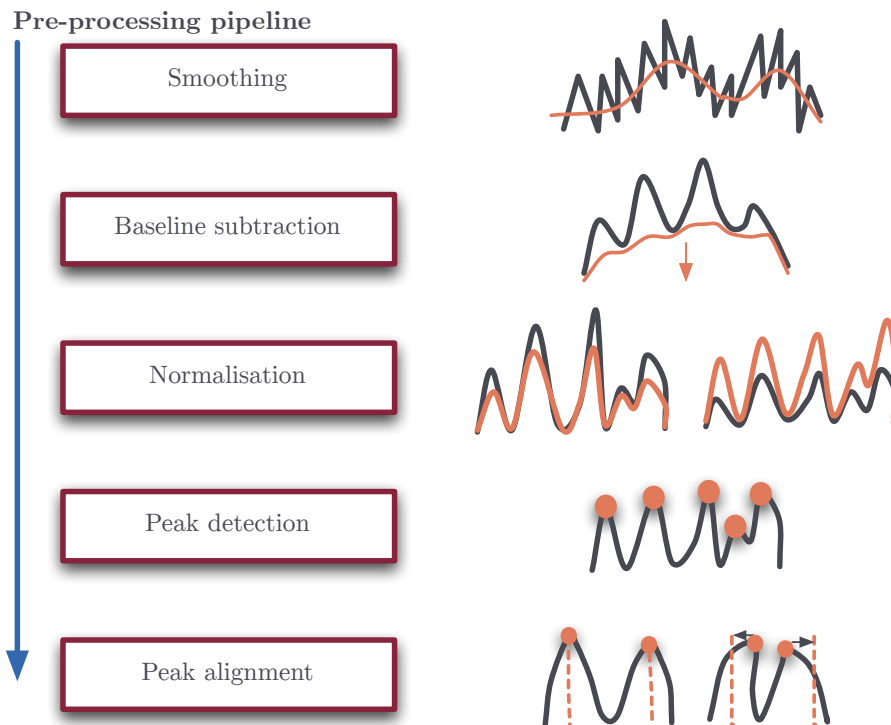


Figure 1: The spectra pre-processing pipeline. The steps, in order, required to successfully pre-process raw proteomic MALDI TOF-MS data.

optimising methods for the baseline subtraction step of pre-processing of the raw spectra.

Baseline subtraction

The non-biological signal to be removed by baseline subtraction is often described as ‘chemical noise’ which predominantly occurs at low mass values and may result from ionised matrix molecules [11]. An example of a MALDI TOF-MS spectrum, a baseline estimate and the resulting baseline subtracted spectrum are shown in Figure 2. The spectrum in Figure 2 is a from the Fiedler dataset which is outlined in the ‘Data used’ section in Methods. The pre-processing applied prior to the baseline subtraction involved taking the square root of the spectrum intensities (for variance stabilisation) and performing the first pre-processing step in smoothing using the Savitzky-Golay method with a half window size of 50 [12].

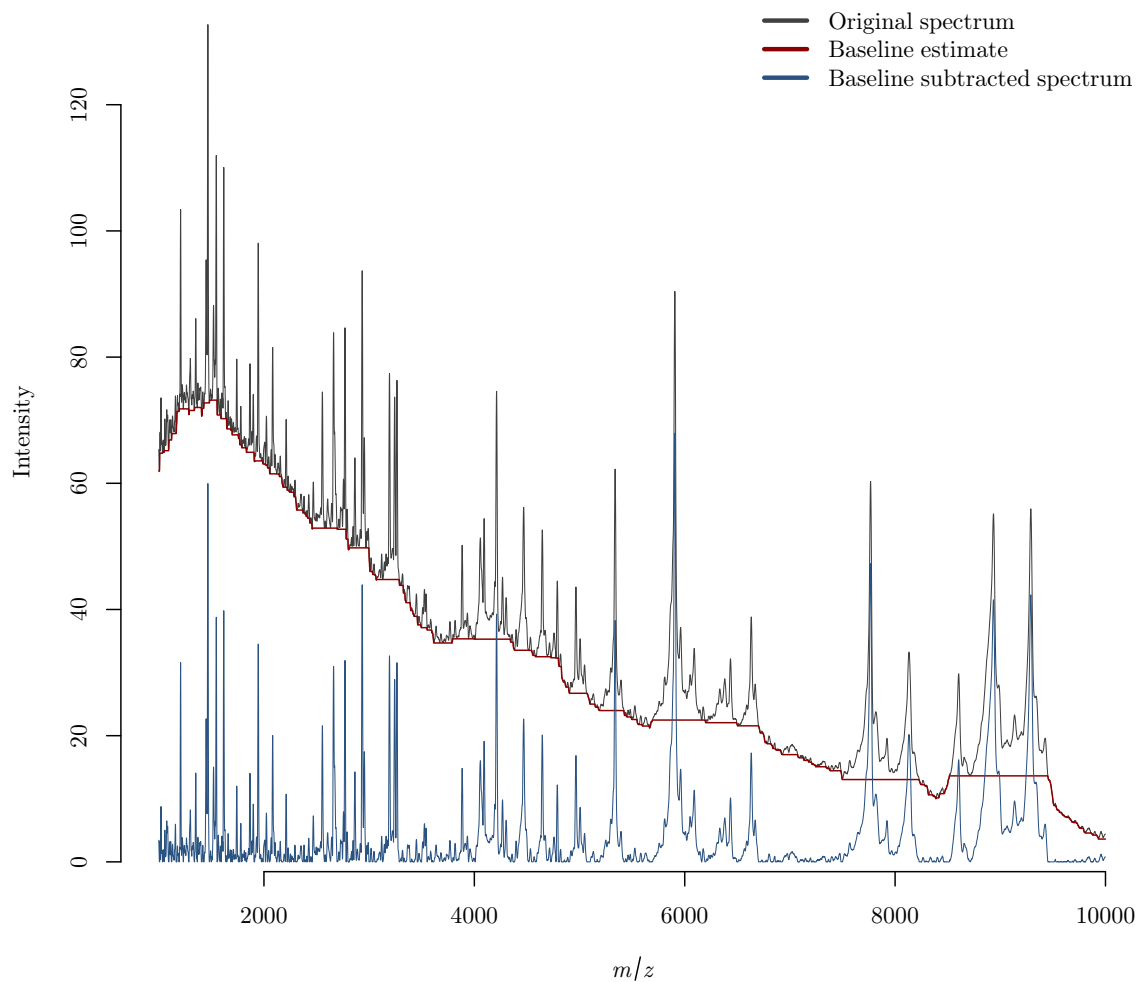


Figure 2: Baseline subtraction of a proteomic MALDI-TOF mass spectrum. A spectrum from the Fiedler dataset: see ‘Data used’ in the Methods section. The square root of the spectrum intensities was taken as a variance stabilisation measure and smoothing using Savitzky-Golay (half window size of 50) was applied prior to baseline subtraction.

The baseline subtraction method discussed in the present paper utilises the *top-hat* operator, which is an operator defined in mathematical morphology. Mathematical morphology was originally proposed for two-dimensional image analysis then further developed for image processing of microarray data images [13, 14]. It has since been applied to MS data [7, 15, 16, 17, 18, 19], and we describe the theory that is largely ignored when applied naively. The mathematical morphology definitions of an erosion, dilation, opening and top-hat provided below allow us to extend the current use of mathematical morphology in MS baseline subtraction.

The top-hat operator has some properties, i.e. it is a non-parametric and non-linear filter, which make it desirable for baseline subtraction. In particular, this suits the non-biological signal in MS spectra which may not follow a known functional form. Furthermore, the top-hat operator is computationally inexpensive compared with standard functional filters that require estimates of model parameters.

Other algorithmic methods of baseline subtraction such as the sensitive nonlinear iterative peak (SNIP) algorithm [20, 21] provide an alternative to the top-hat operator. However, it will be shown in the Methods section that the top-hat operator can importantly be extended, using the mathematical theory underpinning it, for unevenly spaced data.

Standard methods of baseline subtraction estimate local minima (troughs) and fit either local regression (LOESS, Savitzky-Golay) or interpolate (splines) through these points [22]. These methods require careful selection of the window size for detecting troughs, the polynomial order and the span of points for fitting the model, where applicable. Despite using optimised input arguments for these methods, they cannot guarantee a non-negative resultant signal. In fact, padded or removed signal in places of high curvature in the spectra may be produced. This can easily be envisioned by considering two local minimums and an adjacent point to one of the local minimums that lies between both. There is no property that stops the adjacent point lying below an interpolation of the two minimums, especially where there exists a large difference between the values of the local minimums.

Morphological image analysis and theory

The core concepts in mathematical morphology required to apply the top-hat operator are presented below. The definitions of a morphological *structuring element*, *dilation*, *erosion*, *opening* and *top-hat* can also be found in [23, 24, 25, 26].

A structuring element (SE) is a small set that acts on given data or images. For

linear TOF-MS data, a SE is simply a one-dimensional line-segment, or window, passed over the vector of spectral intensities. In the context of morphological image analysis, the SEs used are centred (the median SE value is 0), symmetric (the SE behaves the same either side of the centre) and flat (SE is of the same dimension as the data). Non-flat SEs are not ideal for the current application as they require a known function or weightings to be applied within the sliding window.

Definition 1. *For the sets $X \subset \mathbb{Z}^p$ and $B \subset \mathbb{Z}^p$, $p \in \mathbb{Z}^+$, and the function f defined over X , the erosion of X by B is defined as,*

$$\begin{aligned}\epsilon_B(f)(x) &:= (f \ominus B)(x) \\ &:= \inf_{b \in B} f(x + b),\end{aligned}$$

for each element x in X . The dilation is similarly defined,

$$\begin{aligned}\delta_B(f)(x) &:= (f \oplus B)(x) \\ &:= \sup_{b \in B} f(x + b).\end{aligned}$$

Erosions and dilations can be thought of as rolling minimums and maximums, respectively, over the spectral values. Sometimes the sets X and B in Definition 1 are defined over \mathbb{R}^p [23, 26] but this is rarely implemented for data other than $X \subset \mathbb{Z}^p$ in practice.

Definition 2. *The application of a morphological erosion followed by a morphological dilation to a set X is the morphological opening,*

$$\begin{aligned}\omega_B(f)(x) &:= \delta_B(\epsilon_B(f))(x) \\ &:= ((f \ominus B) \oplus B)(x).\end{aligned}$$

In the context of linear TOF-MS data, a morphological opening is a non-linear estimation of background signal of the one-dimensional spectrum on X . The opening has the desirable property that it is never returns values greater than the observed signal, i.e. $\omega_B \leq f \forall x \in X$.

Definition 3. *The top-hat operator is defined as the removal of the opening from the original signal f ,*

$$\tau_B(f)(x) := f(x) - \omega_B(f)(x).$$

The result of applying the top-hat operator to proteomic TOF-MS is the estimation of the true signal by removing the estimated background signal from f on X . Because of the $\omega_B(f) \leq f$ ($\forall x$) property of morphological openings, the top-hat operator provides a background estimate and removal without risk of creating negative signal, since it is a physical impossibility of the system. Such properties cannot be guaranteed by local regression of local minima.

Example of the top-hat operator

To illustrate the morphological operators that have been defined, we consider a simple example. Let $f = \{a_x\}_{x=1}^{13}$ be a series and define a flat SE, $B = \{b_j\}_{j=1}^5 = \{-2, -1, 0, 1, 2\}$ with

$$f(x) = \begin{cases} a_1 & \text{if } x < 1 \\ a_x & \text{if } x = 1, 2, \dots, 13 \\ a_{13} & \text{if } x > 13, \end{cases}$$

where

$$\{a_x\} = \{ 6 \quad 11 \quad 12 \quad 14 \quad 7 \quad 10 \quad 13 \quad 9 \quad 12 \quad 15 \quad 8 \quad 11 \quad 10 \}.$$

The erosion at $x = 4$ is calculated,

$$\begin{aligned} \epsilon_B(f)(4) &= \inf_{b \in B} f(4 + b) \\ &= \inf \{11, 12, 14, 7, 10\} = 7. \end{aligned}$$

Given the erosions for $x = 2, 3, 4, 5, 6$ are 6, 6, 7, 7, 7, respectively, the morphological opening at $x = 4$ is

$$\begin{aligned} \omega_B(f)(4) &= \delta_B(\epsilon_B(f))(4) \\ &= \sup_{b' \in B} \{\epsilon_B(f)(4 + b')\} \\ &= \sup \left\{ \inf_{b \in B} f(2 + b), \inf_{b \in B} f(3 + b), \inf_{b \in B} f(4 + b), \right. \\ &\quad \left. \inf_{b \in B} f(5 + b), \inf_{b \in B} f(6 + b) \right\} \\ &= \sup \{6, 6, 7, 7, 7\} = 7. \end{aligned}$$

Therefore, the top-hat operator result for $x = 4$ is

$$\tau_B(f)(4) = f(4) - \omega_B(f)(4) = 14 - 7 = 7.$$

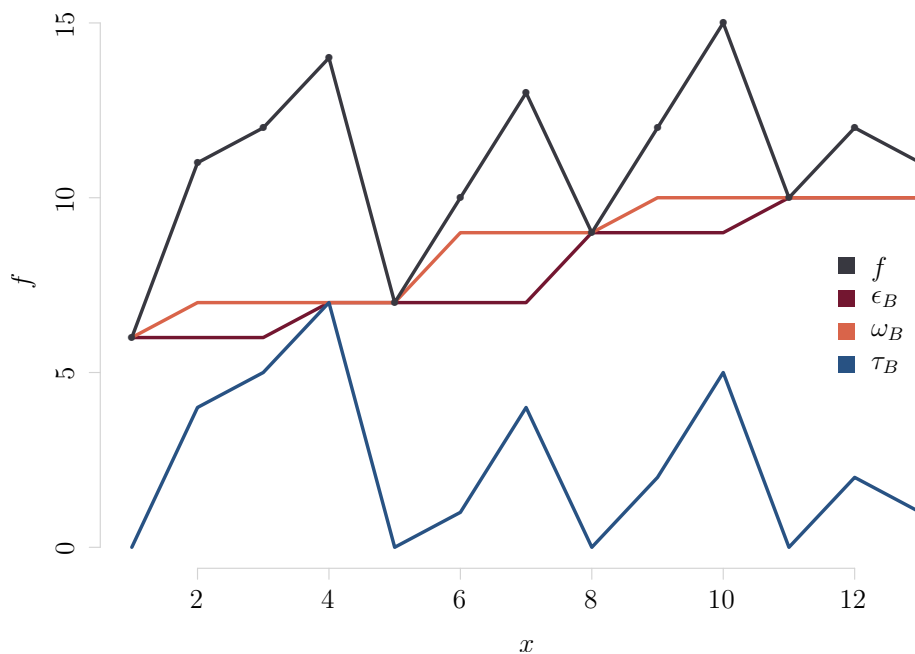


Figure 3: Baseline subtraction on an example spectrum using the top-hat operator (see main text for details): an demonstration of the erosion, opening and top-hat operators (ϵ_B , ω_B and τ_B , respectively) on a set f .

The operations with ϵ_B , ω_B and τ_B using the flat SE, $B = \{-2, -1, 0, 1, 2\}$, on the entire signal $f(x)$ can be observed in Figure 3.

Current application of the top-hat operator to linear TOF-MS

A naive algorithmic application of an erosion to spectral intensities simply requires a traversal of each point, where the minimum value within a window over that point is the resulting erosion. The process is performed similarly for a dilation. However, erosions and dilations can be calculated more efficiently with the line segment algorithm (LSA) [27, 28]. Application of the LSA is mainly seen in medical imaging and analysis [29, 30]. The R package **MALDIquant** and **OpenMS** use this algorithm in their implementation of the top-hat operator.

When applying the top-hat operator to a spectrum, the SE needs to be chosen carefully. In particular, the following need to be considered.

1. If a SE is too large, then it will be too conservative and leave false signal.
2. If a SE is too small, it will result in under-cut peaks and remove valid signal.
3. The mean peak width increases further along the m/z -axis [31]. The baseline subtraction needs to be performed in a piecewise manner, otherwise the above issues 1 and 2 will occur.

Despite the simplicity of the top-hat operator compared to functional alternatives, piecewise baseline subtraction is still required. In fact, piecewise baseline subtraction should be applied for any method that implicitly assumes peak width remains constant, such as local regression, interpolating splines or the SNIP algorithm.

The SE size used for the top-hat operator needs to be of equivalent window size to each spectrum's peak widths, or greater, to ensure the top-hat operator does not 'undercut' peak intensities. The piecewise baseline subtraction involves determining subsections of the m/z -axis, where fixed SE widths (in the number of m/z points) in each section are appropriate, or the equivalent input arguments for other baseline methods. Smaller SEs will be chosen corresponding to lower m/z values and larger SEs will be used corresponding to larger m/z values.

Figure 4(a) illustrates a spectrum from the Fiedler data separated into four roughly equal segments based on the number of intensity values. When applying the top-hat operator, the SE size is a constant number of intensity values within each piecewise section of the axis. The SE sizes selected in Figure 4(a) were made by visual inspection and trial-and-error. Figure 4(b) depicts the same spectrum as Figure 4(a) but the x -axis is in terms of m/z location. On this m/z -axis, the SE size increases along the m/z -axis within each piecewise segment simply by virtue of the distances between m/z points increasing, even though the same window size is being used in terms of the number of intensity values. However, the increasing coverage in m/z units across the m/z -axis is not proportional to the increase in peak widths. Figures 4(a) and (b) demonstrate that there is not a constant number of intensity values for the SE across the entire m/z -axis that could avoid conservative baseline estimates (1) or under-cut peaks (2) or even both.

Improving baseline subtraction

Prior to pre-processing MALDI TOF-MS data, a log or square root transformation of the intensity axis is usually performed as a variance stabilisation measure but no such transformation is made to the m/z -axis. If an appropriate m/z transformation could

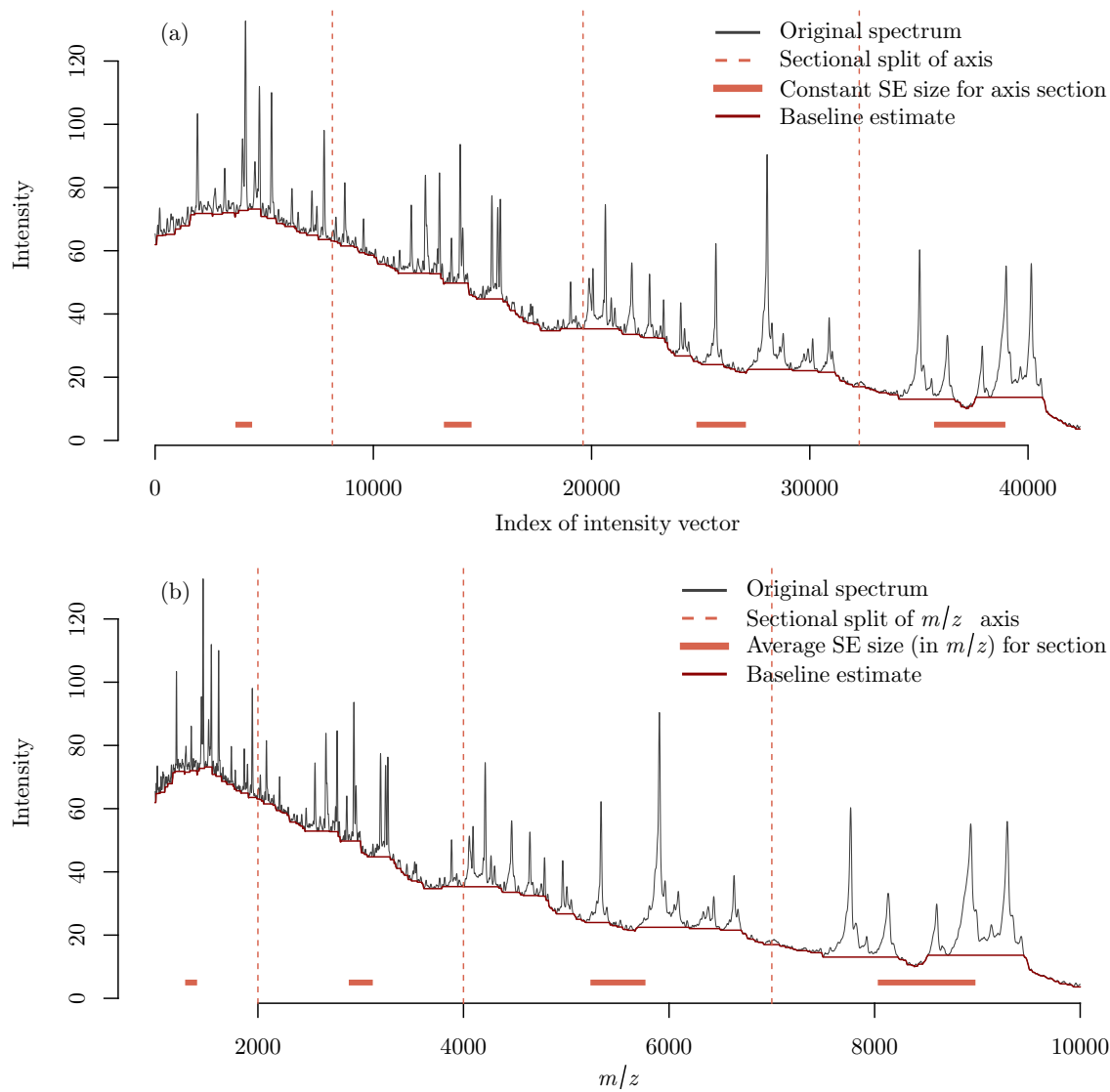


Figure 4: Piecewise baseline subtraction of a proteomic MALDI-TOF mass spectrum from the Fiedler dataset using the top-hat operator. Because the SE is provided as a number of m/z points and m/z values increase in distance along the axis, the SE size is not constant within sub-intervals of the m/z -axis.

be made however, piecewise pre-processing of the spectra for the baseline subtraction step (and potentially for other pre-processing steps) could be avoided. Additionally, the default arguments such as window size in software to perform baseline subtraction are statically defined. Uninformed default arguments such as these are highly likely to need modification for successful baseline subtraction, as spectra attributes vary from one experiment to another. Dynamic default arguments that are informed by the data would be an advantage in saving both user time and minimising user error.

Methods

A pipeline to achieve automated baseline subtraction

The pipeline shown in Figure 5 can be employed to automate the baseline subtraction step. The first step of the pipeline requires a suitable transformation of the m/z -axis. If such a transformation of the m/z -axis can be made, a piecewise approach is not required as a constant-sized SE can be used over the entire spectrum. A log-type transform that expands the low m/z values and contracts the high m/z values is required. Once a suitable transformation is found, a top-hat operator defined over non-evenly spaced real values (i.e. $X \subset \mathbb{R}^p$, as opposed to integer values) can be used at the baseline subtraction step. The implementation requires a minimum and a maximum sliding window algorithm for unevenly spaced data which means the LSA cannot be used. Naive algorithms are available; however, here we present a novel sliding window algorithm that we show outperforms naive sliding window algorithms by avoiding repeated minimum (or maximum) calculations for common points in successive sliding windows. However, a SE size does need to be selected. This can be implemented by firstly estimating peak widths, then selecting a SE size that covers a sufficient proportion of the estimated peak widths. The process of estimating peak widths can be automated without user input and our recommend approach is presented here. The final step in the baseline subtraction pipeline is simply the (reverse) transformation back to the original m/z scale.

The new pipeline to perform baseline subtraction of MALDI TOF-MS data presented in Figure 5 has two major advantages when compared to standard methods.

- Firstly, the pipeline automates the baseline subtraction step, that is otherwise conducted in a piecewise manner. This eliminates the need for user input and time-consuming calibration by observation. Automation of the baseline subtraction step also minimises the potential for user error and the time required

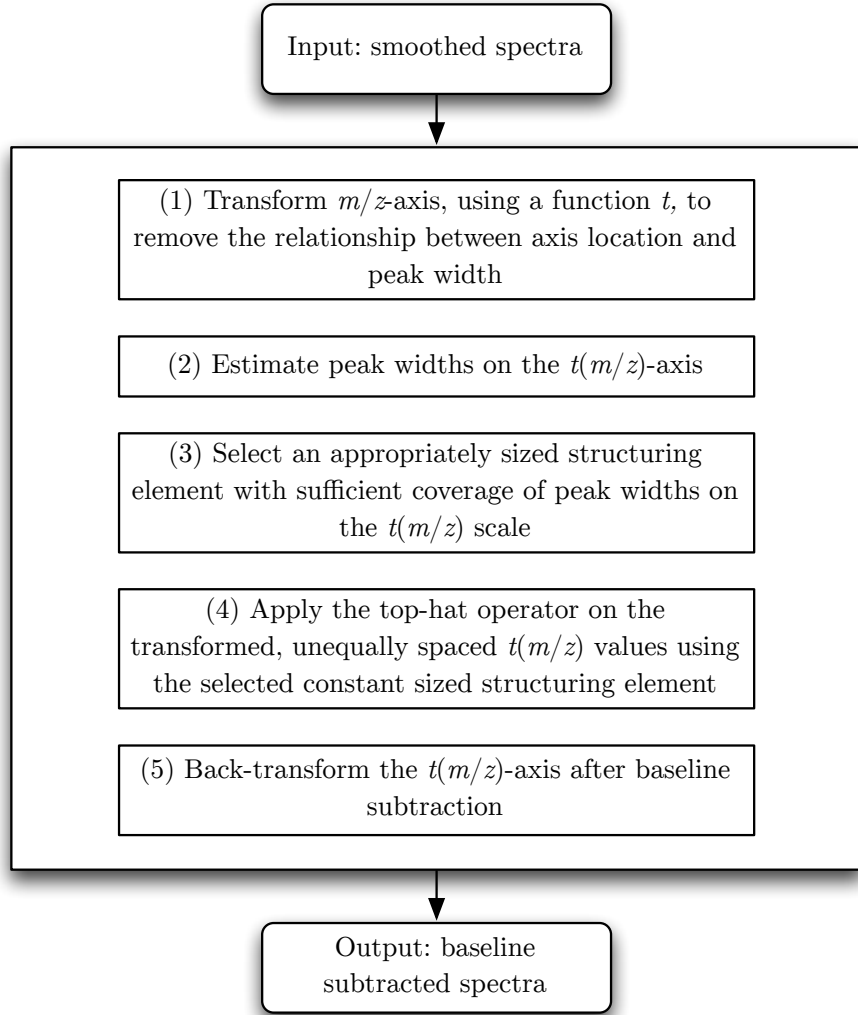


Figure 5: The proposed baseline subtraction pipeline: five steps for automated baseline subtraction.

to assess the input arguments for optimality.

- Secondly, the novel algorithm that is computationally less expensive than a naive minimum or maximum sliding window algorithm to perform the top-hat operation on unevenly spaced data, presented here, further minimises the computational time burden of baseline subtraction.

Fields of application outside of bioinformatics that encounter unevenly spaced data are also likely to find this algorithm useful in practice. Other names for unevenly spaced data include unevenly sampled, non-equispaced, non-uniform, inhomogeneous, irregularly sampled or non-synchronous data. Such data occur in various fields including, but not limited to; financial time-series, geologic time-series, astrophysics and medical imaging [32, 33, 34, 35, 36, 37, 38]. Analysis and processing of unevenly spaced data is an ongoing field of research, as most methods for analysis assume equally spaced data.

Data used

Six proteomic MS datasets from previously published studies were used to validate the methods presented.

Fiedler data: Urine samples were taken from 10 healthy women and 10 healthy men and peptides were separated using magnetic beads (fractionation). The fractionated samples were then subject to MALDI TOF-MS [39]. A subset of the MALDI TOF-MS data is freely available in the R package `MALDIquant` [7] and is the dataset used here. The spectra are observed over the range of values 1,000-10,000 m/z .

Yildiz data: As described in [40], sera were collected from 142 lung cancer patients and 146 healthy controls to find relevant biomarkers. The serum samples were subject to MALDI TOF-MS without magnetic bead separation. The spectra are observed over the range of values 3,000-20,000 m/z .

Wu data: MALDI TOF-MS data were generated from sera, as described in [41, 42], with the aim of differentiating between 47 ovarian and 42 control subjects. The spectra are observed over the range of values 800-3,500 m/z using Reflectron mode which resolve peptide peaks into their isotopomers.

Adam data: Surface-enhanced laser desorption/ionization (SELDI) TOF-MS data from 326 serum samples from subjects classified as prostate cancer, benign hyperplasia or control [43]. While SELDI has been found to be less sensitive than MALDI, samples do not require fractionation before applying MS. The data analysed here are limited to the range 2,000-15,000 m/z as peptide signals beyond this range are sparse.

Taguchi data: The dataset available was first described in [44] but is available as a supplement for [45]. The data are 210 serum-derived MALDI TOF mass spectra from 70 subjects with non-small-cell lung cancer with the aim of predicting response to treatment. The data observed cover the 2,000-70,000 m/z range.

Mantini data: The data in this study were produced using MALDI TOF-MS from purified samples containing equine myoglobin and cytochrome C [46]. A total of 30 spectra are available in the range 5,000-22,000 m/z .

Transformation of the m/z -axis

The proposed pipeline for baseline subtraction requires a suitable transformation of the m/z -axis as the first step. In this section we investigate potential transformations, that will be assessed quantitatively for their suitability.

It has previously been suggested that peak width is roughly proportional to peak location on the TOF-axis [47, 48] and that therefore peak width is proportional to the square of the m/z location. This was not in fact observed for any of the datasets analysed in the present study. Table 1 sets out the shortlist of suitable transformations, t_0 - t_5 , that are considered appropriate for application here.

Table 1: The transforms, t_i , of the m/z -axis trialled to produce a roughly uniform distribution of peak widths across the $t_i(m/z)$ -axis.

Label	Transform
$t_0(x)$	x
$t_1(x)$	$-1000x^{-1}$
$t_2(x)$	$x^{1/4}$
$t_3(x)$	$\ln x$
$t_4(x)$	$-1000(\ln x)^{-1}$
$t_5(x)$	$-1000x^{-1/4}$

To illustrate the role of the transformation, Figure 6 shows a spectrum from the Fiedler dataset on the original m/z -axis (t_0) for transformations t_1 and t_3 . The effect of t_3 , when compared to the original m/z -axis, is an expansion of smaller mass peak widths and the contraction of higher mass peak widths. However, visually it can be seen that higher mass peaks have larger peak widths on average even under the t_3 transformation. The t_1 transformation further shifts low m/z values across the transformed axis and contracts m/z values at the high end of m/z -axis. Potentially, the t_1 transformation creates larger peak widths for smaller m/z values than high m/z values so as to produce peak widths that decrease on average across the transformed axis. The effect of the six transformation functions, t_0 - t_5 , on a spectrum from each of the six datasets is available in Appendix A.

Obtaining approximate peak widths prior to baseline subtraction

Peak widths can be obtained at the peak detection step (step four of pre-processing) but such information is not generally known prior to the second pre-processing step of baseline subtraction. To determine the constant SE size to be passed over the transformed m/z -axis, peak widths need to be estimated. An algorithm to estimate peak widths from the data was created here for this purpose.

The algorithm below to estimate the peak widths within spectra takes the previously smoothed spectra on the transformed m/z -axis as the input and is performed as follows.

- For each spectrum, the lower convex hull of the two-dimensional set of spectrum points is used to determine an approximate baseline for each spectrum.
- The longest segment of the lower convex hull is then halved, with the two sets of points created by this split subject to a new lower convex hull calculation.
- The newly calculated lower convex hull points for the two set of points are then added to the original set of lower convex hull points to improve the approximate baseline calculation.
- This is repeated $r-1$ more times to produce an approximate estimated baseline.
- The approximate baseline is then removed and median intensity is then calculated for the resulting spectrum.

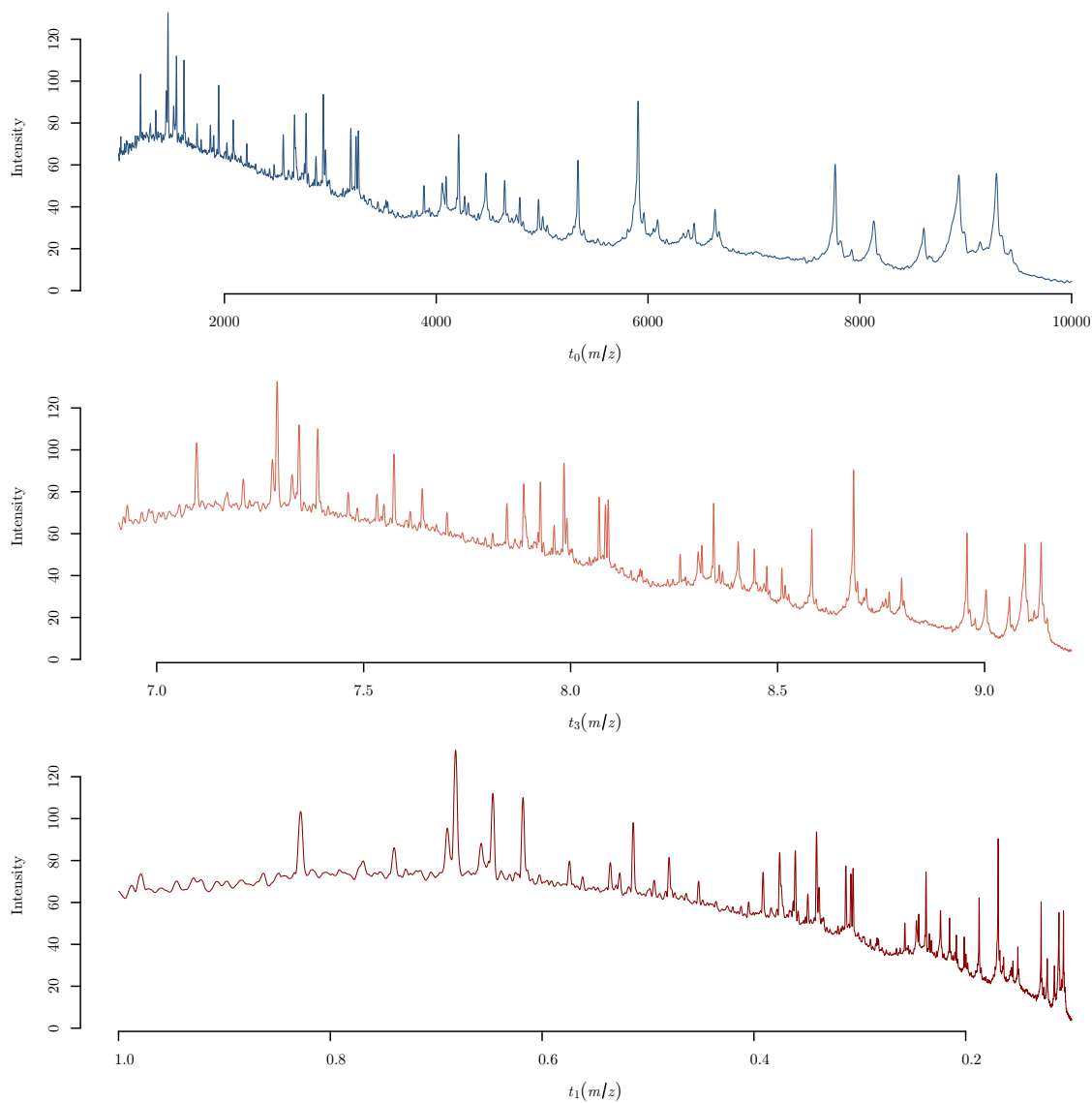


Figure 6: Transformations of the m/z -axis: Three different m/z -axis transformations (see Table 1) for the Fiedler spectrum shown in Figures 2 and 4.

- Intensities above the median value are treated as points along a peak.
- The consecutive points above the median value are the estimated peak widths.

The above algorithm is crude and could not be used for reliable baseline subtraction. However, estimated peak widths are easily extracted using this method and can be used within the proposed automated baseline subtraction pipeline.

A reasonable number of lower convex hull iterations of $r = 5$ produced sensible results on the six datasets used. By specifying a value of r , this method to estimate peak widths is fully automated. It provided enough alterations to the original lower convex hull to satisfactorily remove the residual baseline on concave smoothed spectra while not applying too many alterations so as to create midpoints along the longest segments which create lower convex hulls ending a peak vertices and therefore removing them. However, a missed peak or two per spectrum is not an issue as dozens of peaks are identified per spectrum. Figure 7 depicts this process, on a single spectrum.

The algorithm presented above attempts to automatically find peak widths without user input. We outline this automated procedure in the Methods section as it is not the focus of this paper, and may be substituted with any peak region finding method that requires no user input; such is the modularity of the pipeline shown in Figure 5. There exist other methods to estimate peak widths (regions), such as that found in [20], but they require previous knowledge of likely peak widths and are therefore not a baseline subtraction method that can be automated.

Selecting a SE size and applying the top-hat operator in the transformed space

Point three of Figure 5 requires a choice of SE size. This can be chosen from the estimated peak widths found using the algorithm presented in previous section. The aim is to select a SE of sufficient size to not undercut peaks; such a SE size roughly translates to the maximum of the peak widths. However, there is likely to be a SE size smaller than the maximum estimated peak width but much greater than the minimum estimated peak width that performs optimally. Given a set of estimated peak widths for all spectra in an experiment and a SE size, we define the proportion of peak widths that are estimated to be the SE size or smaller as the estimated peak coverage proportion (EPCP) .

Figure 8 represents the estimated peak widths for the 16 spectra in the Fiedler

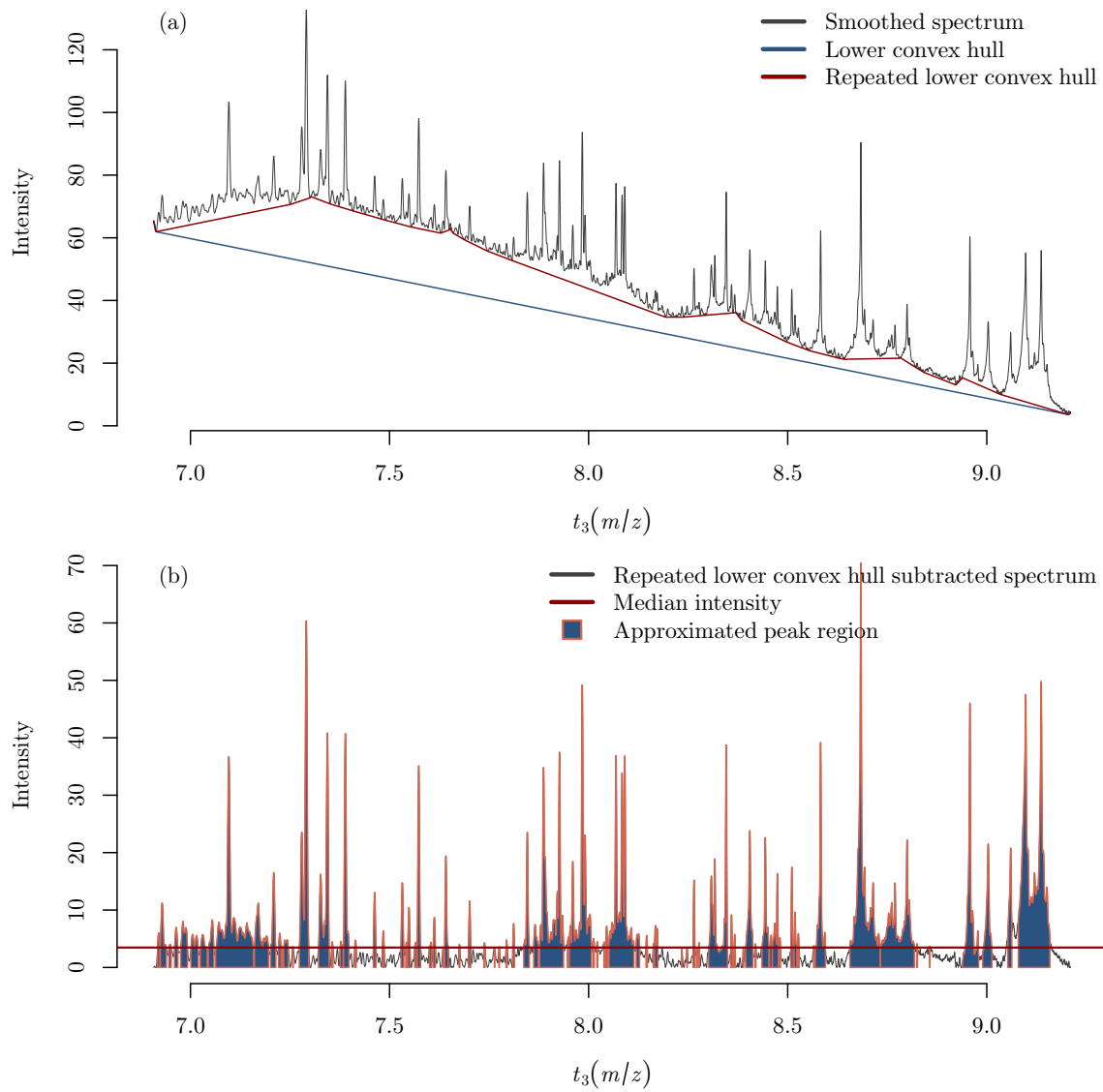


Figure 7: Illustrated algorithm to determine approximate peak widths prior to baseline subtraction. (a) A Fiedler spectrum with the lower convex hull and repeated ($r = 5$) lower convex hulls shown. (b) The lower convex hull subtracted spectrum with median intensity of the spectrum depicted. Points above the median are considered peaks, which have been filled with colour.

dataset on the $t_2(m/z)$ scale, where peak regions are found using a repeated lower convex hull algorithm presented previously.

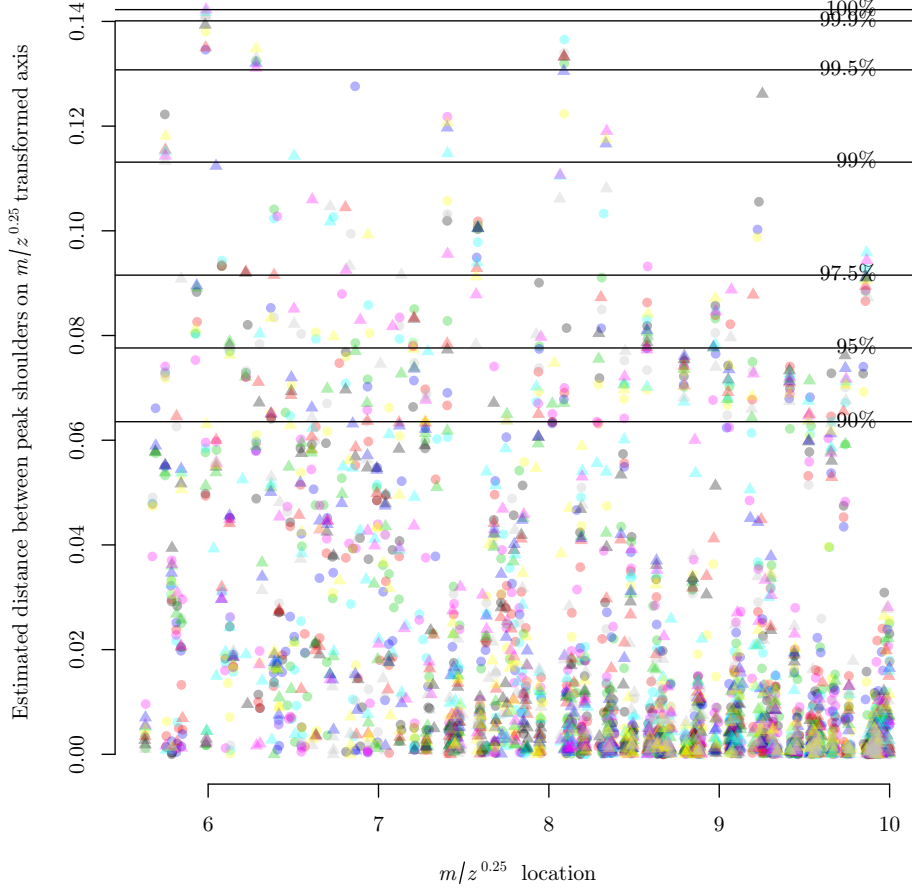


Figure 8: Estimated peak widths on the $t_2(m/z)$ scale. Peak widths in the Fiedler data with the m/z -axis transformed by the quartic root. The horizontal lines denote the proportion of peak widths that lie below it.

We trial different SE sizes corresponding to different EPCP values in the hope an optimal EPCP value for each of the six datasets we utilise can be found. A SE size that fully covers 97.5% of detected estimated peak widths (EPCP of 0.975) for example, could yield optimised baseline subtraction.

Both the EPCP and m/z -axis transformation are variables that are explored in the Results section, to find an empirically optimal combination. Optimality of the

automated baseline subtraction can be assessed by calculating the minimum value of an error metric relative to a gold-standard baseline subtraction, given a set of EPCP values and transformation functions. The metric used to compare the automated baseline subtraction to the gold-standard is outlined in the next section and the modified algorithm to perform top-hat baseline subtraction on the unevenly spaced and transformed m/z -axis is provided in the section after that.

Comparison of proposed methods to the gold-standard

Piecewise, top-hat baseline subtracted spectra were used as the gold-standard baseline subtracted spectra. The SE sizes for each piecewise segment along the m/z -axis were selected using trial-and-error to produce the best baseline subtraction as determined visual inspection. These baseline subtracted, gold-standard spectra were produced prior to the automated baseline subtraction methods being applied.

Mean absolute scaled error (MASE [49]) was selected to be the error metric of the automatically baselined spectra for a given transformation and EPCP, when compared to the gold-standard baseline subtracted spectra. Because the MALDI TOF mass spectra intensities are on arbitrary scales prior to normalisation, it is important to use a metric that is scale free, in order to be able to compare results between spectra from different experiments. MASE also avoids many degeneracy issues of other relative error metrics with zero denominators. Baseline subtracted spectra will have many zero values where no signal is present. Other metrics such as mean squared error (MSE) were considered (which did not change the selection of the optimal transform and EPCP) however the ability to compare the error with other data is not possible and some sort of normalisation or weighting of spectra is required to ensure the MSE, say, of selected spectra do not dominate the result.

Let τ_j^* denote the intensity at x_j of a gold-standard baseline subtracted spectrum $\tau_{B^*}(x_j)$ and τ_j denote an automated baseline subtracted spectrum $\tau_B(x_j)$. The MASE is calculated as

$$\text{MASE} = \text{mean} \left(\left\{ \frac{|\tau_j^* - \tau_j|}{\frac{1}{n-1} \sum_{i=2}^n |\tau_i^* - \tau_{i-1}^*|} \right\}_{j=1,2,\dots,n} \right).$$

For each of the six datasets, there are N spectra to be compared. Let AMASE be

the average MASE value of the N baseline subtracted spectra, then

$$\text{AMASE} = \frac{1}{N} \sum_{\ell=1}^N \text{MASE}_{\ell}.$$

The ‘continuous’ line segment algorithm

A novel algorithm is proposed here that can be applied to the unevenly spaced values of the transformed m/z -axis using a constant SE width. This algorithm, which we name the ‘continuous’ line segment algorithm (CLSA), requires fewer computations per element than current rolling maximum and minimum algorithms on unevenly spaced data [50].

Consider the case where values in X are not evenly spaced, and $X \subset \mathbb{R}$, as opposed to $X = \{1, 2, \dots, n\}$, such as proteomic spectra on a transformed $t(m/z)$ -axis. Figure 9 outlines the CLSA as a rolling minimum algorithm that can be trivially converted to a rolling maximum algorithm by finding the rolling minimum of $-f$ and returning the negative values of the result.

In effect, the CLSA creates m blocks using the θ_i relating to the corresponding x_i :

$$\begin{aligned} \theta_1, \theta_2, \dots, \theta_{b_1} &= 1 && \text{where } x_1, x_2, \dots, x_{b_1} \in [x_1, x_1 + k) \\ \theta_{b_1+1}, \theta_{b_1+2}, \dots, \theta_{b_2} &= 2 && \text{where } x_{b_1+1}, x_{b_1+2}, \dots, x_{b_2} \in [x_1 + k, x_1 + 2k) \\ &\vdots && \\ \theta_{b_{m-1}+1}, \theta_{b_{m-1}+2}, \dots, \theta_{b_m} &= m && \text{where } x_{b_{m-1}+1}, x_{b_{m-1}+2}, \dots, x_{b_m} \in [x_1 + (m-1)k, x_n]. \end{aligned}$$

When the algorithm considers each point x_i for the minimum f in the window spanning $k/2$ either side, it checks whether the most extreme x -values in this window are either in the current block or one block away (these values cannot be further than one block away as block sizes are of length k) to decide on which combination of g and h is required. Note the algorithm is impervious to arbitrarily spaced x_i as long as they are in ascending order. If $\theta_i \neq j$ for any $i = 2, 3, \dots, n-1$; $j = 2, 3, \dots, m-1$ (empty blocks) or $x_{b_{j-1}+1} = x_{b_j}$ for any $j = 2, 3, \dots, m$ (blocks with only one x_i), for example, do not affect the validity of the proposed algorithm.

This algorithm can be seen as a generalised version of the LSA [27, 28] as it works on evenly and unevenly spaced data. An R implementation of this novel CLSA can be found as an R-package using compiled C code at <https://github.com/tystan/clsa>.

Let $k_0 = \frac{k}{2}$ where k is the length of a centred, one-dimensional window. Consider the ordered (ascending) set of transformed m/z points, $X = \{x_1, x_2, \dots, x_n\}$, and the corresponding expression values, f . Furthermore, define $\text{span}(X) = x_n - x_1$ where $\text{span}(X) > k$ and choose the smallest $m \in \mathbb{Z}^+$ so that $mk \geq \text{span}(X)$.

Three vectors taking integer values are required to be created initially,

$$\Theta = [\theta_1, \theta_2, \dots, \theta_n], W^\nabla = [w_1^\nabla, w_2^\nabla, \dots, w_n^\nabla], W^\Delta = [w_1^\Delta, w_2^\Delta, \dots, w_n^\Delta].$$

For $i = 1, 2, \dots, n$, the integer θ_i is calculated as follows,

$$\theta_i = \{j : \text{if } x_1 + (j-1)k \leq x_i < x_1 + jk \text{ for } j \in \{1, 2, \dots, m\}\};$$

w_i^∇ is the index corresponding to x_i satisfying the inequality,

$$x_{w_i^\nabla-1} < x_i - k_0 \leq x_{w_i^\nabla};$$

and w_i^Δ is the index corresponding to x_i satisfying the inequality,

$$x_{w_i^\Delta} \leq x_i + k_0 < x_{w_i^\Delta+1}.$$

The rolling minimum at x_i can be calculated as,

$$r_{\min}(f(x_i)) = \begin{cases} g(x_{w_i^\Delta}) & \text{if } \theta_{w_i^\nabla} = \theta_{w_i^\Delta+1} \\ h(x_{w_i^\nabla}) & \text{if } \theta_{w_i^\nabla-1} = \theta_{w_i^\Delta} \\ \min\{g(x_{w_i^\Delta}), h(x_{w_i^\nabla})\} & \text{otherwise,} \end{cases}$$

where

$$\begin{aligned} g(x_i) &= \begin{cases} f(x_i) & \text{if } \theta_{i-1} < \theta_i \text{ (define } \theta_0 = 0) \\ \min\{g(x_{i-1}), f(x_i)\} & \text{otherwise; and} \end{cases} \\ h(x_i) &= \begin{cases} f(x_i) & \text{if } \theta_i < \theta_{i+1} \text{ (define } \theta_{n+1} = m+1) \\ \min\{f(x_i), h(x_{i+1})\} & \text{otherwise.} \end{cases} \end{aligned}$$

Figure 9: The continuous line segment algorithm (CLSA).

A demonstration of why the creation of blocks the size of the SE and accessing cumulative values half an SE length away allows the calculation of rolling minimums is shown in [27]. Examples to demonstrate the mechanics of the CLSA algorithm are presented in Appendix B.

Results

Presented in this paper is a pipeline to automate the baseline subtraction step in proteomic TOF-MS pre-processing. The pipeline consists of transforming the m/z -axis, then finding an appropriate SE size via an automated peak width estimation algorithm on the transformed scale, applying a novel algorithm to perform the top-hat baseline subtraction, then finally, baseline subtracted spectra are returned by back-transforming the data to the m/z scale.

There remain two elements of the pipeline to be assessed. Firstly, for the pipeline to be fully automated, an optimal combination of EPCP value and transformation need to be found. In the next section we perform a grid search over EPCP values of 0.8, 0.85, 0.9, 0.95, 0.98, 0.99, 1 and transformations $t_0, t_1, t_2, t_3, t_4, t_5$ to find which combination provides the closest baseline subtracted signal to the gold-standard. Given sufficient similarity to the gold-standard is achieved, it is hoped that a consensus over all datasets, in their varying attributes, of the optimal combination of EPCP value and transformation can be found. If a consensus is indeed found, the pipeline is likely to be applicable to other proteomic TOF-MS datasets.

A theoretical and empirical assessment of the efficiency of the CLSA in comparison to naive rolling window algorithms then follows. The theoretical efficiency is discussed with respect to the number of operations required over all the elements input into the CLSA. By performing the top-hat operation on the six proteomic TOF-MS datasets and simulated datasets of varying sizes, the computational time required for the CLSA versus the naive algorithm provides an empirical assessment of their relative efficiencies.

Comparison of piecewise and transformed axis baseline subtraction

Figure 10 and Table 2 present the AMASE values on the six datasets. For each dataset, Figure 10 displays a grid of the input arguments which are the m/z -axis

transformation and the EPCP. The colour of the blocks at the intersection of these combinations depict the AMASE value obtained. Darker blocks indicate smaller, and thus preferred, AMASE values. Table 2 is simply a tabular presentation of the results shown in Figure 10.

No single transformation or EPCP was optimal. However, EPCP between 0.95 and 0.99 provided the optimal AMASE value for all datasets suggesting the peak width estimation process is relatively stable. On the Fiedler, Yildiz, Taguchi and Mantini datasets, the null transformation which implicitly implies a constant peak width across the m/z -axis is not valid as AMASE values are notably higher than for the remaining transformations. The transformations t_2 , t_3 , t_4 and t_5 produced the best results. It should be noted that the transformations t_3 , t_4 and t_5 produced very similar AMASE values. With the exception of the Yildiz dataset, using these transformations with an EPCP of 0.95 produced sensible results.

Figure 11 demonstrates the baseline estimates using the gold-standard piecewise top-hat operator, the AMASE optimal transformation and EPCP (t_3 , 0.98) and a non-optimal combination of transformation and EPCP (t_4 , 0.95) that was suitable on all but the Yildiz data. The optimal AMASE transformation and EPCP combination (t_3 , 0.98) shows very little difference from the gold-standard baseline estimate.

Because the gold-standard baseline estimate is subject to expert input and opinion, the differences seen in the gold-standard and the optimal AMASE baseline estimate are not of concern as both look sensible. The non-optimal baseline estimate produces a reasonable automated baseline subtraction, however, it can be seen that this estimate does undercut the peaks especially at high m/z -values.

With respect to the AMASE values, the spectra with fewer peaks generally had larger AMASE values; this is a function of the normalising constant for each spectra, $\frac{1}{n-1} \sum_{i=2}^n |\tau_i^* - \tau_{i-1}^*|_{j=1,2,\dots,n}$, as fewer peaks will generally imply less relative change in signal.

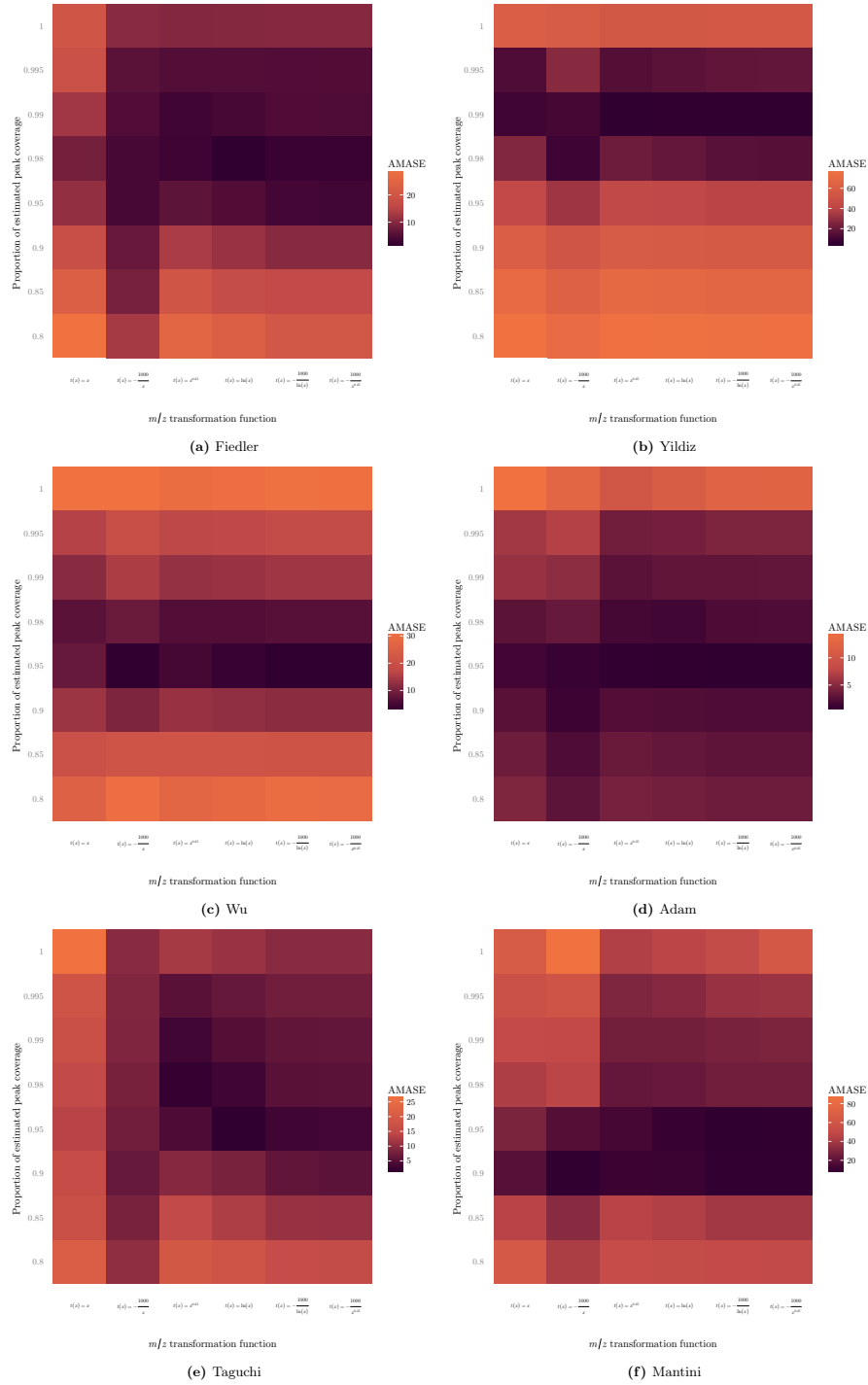


Figure 10: Average mean absolute scaled error (AMASE) heatmaps: AMASE values for each of the six datasets under different combinations of m/z transformation and estimated peak coverage proportion (EPCP).

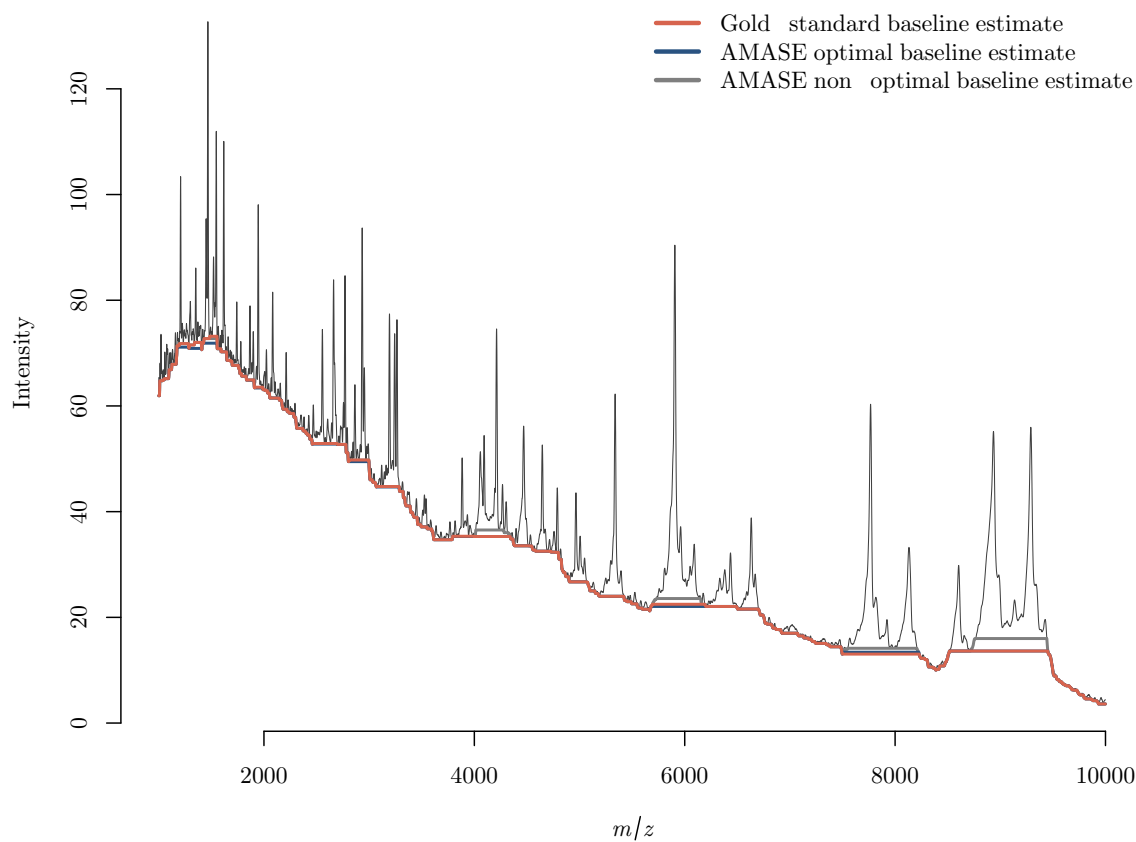


Figure 11: Optimised automated baseline estimate (blue) in comparison to the gold-standard (orange) piecewise baseline estimate for the Fiedler spectrum: the optimal transformation and EPCP were t_3 and 0.98, respectively; the non-optimal combination (grey) of transformation and EPCP shown is t_4 and 0.95, respectively.

Table 2: Average mean absolute scaled error (AMASE) when using structuring element (SE) sizes corresponding to different estimated peak coverage proportions (EPCPs) for each of the selected short-listed transformations on each of the six datasets.

EPCP	$t_0(x) = x$	$t_1(x) = \frac{-1000}{x}$	$t_2(x) = x^{1/4}$	$t_3(x) = \ln x$	$t_4(x) = \frac{-1000}{\ln x}$	$t_5(x) = \frac{-1000}{x^{1/4}}$
Fiedler						
1	19.8	10.1	9.5	9.7	9.9	9.9
0.995	18.4	5.5	4.8	4.7	4.7	4.6
0.99	12.4	4.7	3.0	3.7	4.5	4.5
0.98	8.2	3.7	2.8	1.3	2.3	2.3
0.95	11.0	3.7	5.8	4.6	3.3	3.1
0.9	17.8	6.9	13.2	11.7	9.9	10.0
0.85	23.2	8.3	19.8	16.9	16.1	16.2
0.8	29.4	12.9	25.5	22.9	20.2	20.2
Yildiz						
1	59.7	58.3	54.9	54.6	55.6	55.7
0.995	12.9	27.4	14.3	15.8	17.5	18.0
0.99	9.3	10.7	5.0	5.2	5.5	5.5
0.98	25.7	8.6	20.9	18.4	15.4	14.9
0.95	42.3	32.6	41.6	41.0	39.7	39.5
0.9	59.7	52.7	58.4	57.7	56.7	56.6
0.85	70.4	62.9	69.2	68.1	67.2	67.0
0.8	77.1	71.5	76.6	75.9	75.0	74.9
Wu						
1	30.9	30.8	29.2	29.7	30.4	30.2
0.995	16.1	19.5	16.9	17.2	18.0	17.9
0.99	12.0	15.2	13.0	13.5	14.1	14.0
0.98	7.4	9.0	6.7	6.7	7.2	7.1
0.95	8.8	3.4	5.4	4.1	3.3	3.3
0.9	13.8	10.9	13.3	12.8	12.3	12.3
0.85	20.1	20.9	20.9	21.1	20.7	20.7
0.8	25.4	29.7	27.0	27.7	28.8	28.7
Adam						
1	14.8	12.7	10.1	11.2	12.3	12.4
0.995	6.2	7.1	3.8	3.9	4.5	4.5
0.99	5.7	5.2	2.7	3.0	3.1	3.1
0.98	2.8	3.4	1.6	1.5	2.1	2.1
0.95	1.5	1.0	0.7	0.6	0.6	0.5
0.9	2.6	1.2	2.3	2.2	2.1	2.1
0.85	3.7	2.2	3.4	3.2	2.9	2.9
0.8	4.6	2.9	4.1	3.9	3.7	3.7
Taguchi						
1	26.5	9.6	12.0	11.0	9.6	9.7
0.995	17.7	8.9	5.5	6.5	7.5	7.5
0.99	16.4	8.8	3.5	5.1	6.1	6.2
0.98	14.8	8.1	2.2	3.3	5.5	5.6
0.95	13.6	7.6	4.5	1.8	3.3	3.6
0.9	15.7	6.7	9.4	8.2	6.2	5.8
0.85	16.5	8.2	14.7	12.7	10.8	10.7
0.8	20.7	10.2	18.9	17.7	15.4	14.9
Mantini						
1	66.1	86.4	42.8	45.6	50.6	63.0
0.995	55.1	58.2	29.6	31.8	36.3	37.2
0.99	49.3	48.3	26.2	26.2	27.9	28.9
0.98	42.1	45.9	22.3	23.4	25.8	25.9
0.95	28.8	18.1	13.9	9.5	7.6	7.9
0.9	18.8	8.1	11.0	10.6	8.5	8.5
0.85	45.3	32.4	45.0	42.9	38.9	38.9
0.8	63.1	41.6	51.2	49.7	49.2	49.0

Efficiency of the CLSA compared to the naive rolling window

The naive rolling minimum algorithm consists of the linear-time process of finding the indexes of points at the upper and lower edges of the sliding window for each element, by incrementing the edge indexes from the previous element when required. Using a_k as the average number of data points in the sliding window of size k , the computational cost of finding the minimum value in the window requires approximately $a_k - 1$ comparisons per element. This is because each element requires, on average, a minimum or maximum comparison of all the data points in the window except one: the first data point does not require a comparison. The resulting computational complexity is $\mathcal{O}(a_k n)$ for the naive algorithm, which is dependent on the size of the sliding window and the number of elements in X .

Like the LSA, the CLSA is a linear-time algorithm irrespective of the window size, k . For the CLSA, a linear-time progression through the n elements is required to assign integers of the Θ -vector, as each element is an integer equal to or greater than that which precedes it. The linear-time process of finding the W^∇ and W^Δ indexes at the lower and upper edges of the sliding window, respectively, for each element is similar to that required in the naive algorithm. One linear-time sweep forward and one linear-time sweep back on the data is required to create g and h . A final sweep of the created vectors W^∇ , W^Δ , Θ , g and h is required to compute the r_{\min} values. Each $r_{\min}(f(x_i))$ calculation requires the tests $\theta_{w_i^\nabla} = \theta_{w_i^\Delta+1}$, $\theta_{w_i^\nabla-1} = \theta_{w_i^\Delta}$ or $\min\{g(x_i), h(x_i)\}$. It can therefore be deduced the CLSA is $\mathcal{O}(n)$ complexity, requiring a series of linear-time operations, importantly independent of the length of the sliding window, k .

Given the MS application, $a_k - 1$ operations per element in the naive algorithm would be much larger than the constant number of operations required per element for the CLSA and efficiency strongly favours the CLSA. It should be pointed out that the CLSA requires extra memory availability beyond the iterative algorithm for the creation of the vectors W^∇ , W^Δ , Θ , g and h . Another computational advantage of the CLSA is that by using the minimum of the two temporary vectors g and h as opposed to the minimum of a non-constant number of data points for each $x_i \in X$, vectorised programming can be utilised instead of loops. This is of significant advantage in programming languages that are interpreted such as R.

Using the `clsa` package, the CLSA and naive sliding window algorithms were compared for computational time to calculate the top-hat on real and simulated data. The computations were performed on a 21.5" iMac (late 2013 model, 2.7GHz Intel Core i5, 8GB 1600MHz DDR3 memory, OS X 10.10.2). To optimise speed, the cal-

culations requiring iterative looping were performed using compiled C code for both the CLSA and naive algorithms. The code to run the test of computational running time on the simulated data is provided in Appendix C.

The CLSA and naive sliding window algorithms were applied to perform top-hat baseline estimation to the six datasets used in this paper and the results are shown in Table 3. The CLSA resulted in a reduction of the required computational time by a factor of at least 4. The advantage in speed of the CLSA had greater improvement for the datasets with a greater number of m/z values. The biggest relative improvement was by a factor exceeding 50 for the largest dataset in terms of m/z values per spectra on the Yildiz spectra.

Table 3: Computational time to perform top-hat baseline subtraction in the transformed space using the naive and CLSA algorithms on the six datasets under study.

Data	Number of spectra	Number of m/z values	Computational time (sec)	
			Naive algorithm	CLSA
Fiedler	16	42388	7.7	0.2
Yildiz	264	75958	312.6	5.5
Wu	89	91378	34.1	1.7
Adam	326	8461	3.0	0.7
Taguchi	210	19234	18.0	0.9
Mantini	30	32967	6.7	0.2

Table 4 displays the computational times of top-hat baseline estimation using the CLSA and naive algorithms for varying datasets and SE sizes. The simulated data consisted of 20 randomly generated spectra with x_i and f_i for $i = 1, 2, \dots, n$. These values were independently and randomly generated, where the signal locations $x_i \sim \text{Beta}(1, 3)$ mimic a higher density of points at the low end of the spectra and $f_i \sim \chi_{10}^2$ mimic the positive signals in spectra. MALDI TOF-MS data can have in excess of tens of thousands of m/z values, hence, values of $n = 10^4, 2 \times 10^4, \dots, 10^5$ were used. Varying window sizes were tested, ranging in width from 0.5% to 20% of the x -axis domain. i.e., 0.5% corresponds to a window size of 0.005 passed over the domain $[0, 1]$.

The CLSA was faster than the naive algorithm in every scenario as shown in Table 4. As expected, the computational time was constant for the CLSA irrespective of the window size for a fixed number of points (number of transformed m/z values). The difference in computational time between the two algorithms was reasonably small for

Table 4: Computational time in seconds to perform top-hat baseline subtraction in the transformed space using the naive and CLSA algorithms on synthetic data for varying data assumptions and SE sizes.

Number of points $n (\times 10^4)$	Naive						CLSA					
	Window size (% of x -axis)						Window size (% of x -axis)					
	0.5	1	2	5	10	20	0.5	1	2	5	10	20
1	0.1	0.1	0.3	0.6	1.2	2.3	0.0	0.0	0.2	0.0	0.0	0.2
2	0.3	0.5	1.0	2.5	4.9	9.2	0.1	0.1	0.1	0.1	0.1	0.1
3	0.6	1.2	2.3	5.7	11.0	20.6	0.3	0.1	0.1	0.2	0.1	0.1
4	1.1	2.1	4.2	10.2	19.6	36.5	0.1	0.2	0.1	0.1	0.3	0.1
5	1.7	3.3	6.5	15.8	30.5	56.8	0.2	0.3	0.2	0.2	0.3	0.2
6	2.4	4.7	9.3	22.7	44.0	82.0	0.2	0.2	0.3	0.2	0.2	0.3
7	3.2	6.4	12.6	30.9	59.7	111.7	0.2	0.2	0.4	0.2	0.2	0.4
8	4.2	8.4	16.6	40.3	78.3	146.4	0.4	0.3	0.3	0.4	0.3	0.3
9	5.4	10.6	21.0	51.1	98.8	185.3	0.4	0.3	0.3	0.4	0.3	0.3
10	6.6	13.0	25.9	63.1	121.8	228.7	0.3	0.4	0.3	0.3	0.5	0.3

small datasets and small SE sizes. However, for a typical number of m/z points seen in practice, say 50,000, and a moderate window size that on average encapsulates 5,000 points (1% of x -axis), the CLSA provides an order of magnitude increase in speed.

Discussion

The current gold-standard in baseline subtraction is a piecewise approach that is performed manually, that is, by inspection. Piecewise baseline subtraction is typically performed because, as we have consistently observed with the datasets analysed in this paper, the properties of the spectra do not remain constant over their domain. In particular, a spectrum’s peak width increases with increasing m/z -values. We have proposed a new baseline subtraction pipeline be adopted for the correction of mass proteomic spectra data which avoids both the manual user input and the piecewise-subtraction aspect of existing methods. Our new pipeline is based on the premise that a suitable transformation of the m/z -axis can be found which removes the relationship between peak width and peak location.

As part of the new pipeline, we propose a method to create data-based, and therefore data specific, peak-width estimates from smoothed spectra. Even if this step is not used to automate baseline subtraction, it provides an initial sensible SE size that adapts to each individual dataset. Our generalised version of the LSA is also

presented in the paper, which we call CLSA. CLSA can be applied to unevenly or evenly spaced data and is not limited in its application to proteomic MS data. Should a transformation be known to create peak widths independent of m/z -location in proteomic MS data, an efficient and effective baseline subtraction can be performed using the top-hat operator with a CLSA implementation. A major contribution to note is that we have demonstrated CLSA far outperforms the naive rolling minimum algorithm in required computational time by an order of magnitude or more on numerous datasets of real-world complexity.

The transformed and constant-sized window approach may suffer from a slight but largely unnoticeable reduction in sensitivity in comparison. The trade-off between exactness of the piecewise approach and the speed of the automated transformation and continuous approach may be a consideration, especially if a known relationship exists between the peak width and peak location.

Availability of supporting data

Fiedler data: A subset of the MALDI TOF-MS data generated by the study [39] is available in the publicly available R package: `MALDIquant` [7].

Yildiz data: Available at

<http://www.vicc.org/biostatistics/serum/JT02007.htm>.

Wu data: Previously available at

<http://bioinformatics.med.yale.edu/MSDATA>.

Adam data: Data was obtained on request from the authors of [43]. However some Eastern Virginia Medical School data is available at

<http://edrn.nci.nih.gov/science-data>.

Taguchi data: Available at

<http://www.vicc.org/biostatistics/download/WSData.zip>.

Mantini data: Available at

<http://www.biomedcentral.com/content/supplementary/1471-2105-8-101-S2.zip>.

Competing interests

The authors declare that they have no competing interests.

Author’s contributions

TS and PS developed the statistical and analytical methods. CB and PS provided guidance on the analysis of proteomic data. TS developed the code and implementation. All authors contributed to the writing of the manuscript.

Acknowledgements

Thank you to the creators and custodians of the publicly available data used in this manuscript. TS’s PhD research was funded by an Australian Postgraduate Award scholarship.

References

- [1] J Albrethsen. Reproducibility in protein profiling by MALDI-TOF mass spectrometry. *Clinical Chemistry*, 53(5):852–858, 2007.
- [2] V Kulasingam and E P Diamandis. Strategies for discovering novel cancer biomarkers through utilization of emerging technologies. *Nature Clinical Practice Oncology*, 5(10):588–599, 2008.
- [3] G L Hortin. The MALDI-TOF mass spectrometric view of the plasma proteome and peptidome. *Clinical Chemistry*, 52(7):1223–1237, 2006.
- [4] A Croxatto, G Prod’hom, and G Greub. Applications of maldi-tof mass spectrometry in clinical diagnostic microbiology. *FEMS Microbiology Reviews*, 36(2):380–407, 2012.
- [5] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014.
- [6] R C Gentleman, V J Carey, D M Bates, and others. Bioconductor: Open software development for computational biology and bioinformatics. *Genome Biology*, 5:R80, 2004.
- [7] S Gibb and K Strimmer. MALDIquant: a versatile R package for the analysis of mass spectrometry data. *Bioinformatics*, 28(17):2270–2271, 2012.

- [8] X Li. *PROcess: Ciphergen SELDI-TOF Processing*, 2005. R package version 1.42.0.
- [9] C A Smith, E J Want, G O’Maille, R Abagyan, and G Siuzdak. XCMS: Processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching and identification. *Analytical Chemistry*, 78:779–787, 2006.
- [10] T E Stanford. *Statistical analysis of proteomic mass spectrometry data for the identification of biomarkers and disease diagnosis*. PhD thesis, School of Mathematical Sciences, The University of Adelaide, 2015.
- [11] G L Glish and R W Vachet. The basics of mass spectrometry in the twentyfirst century. *Nature Reviews Drug Discovery*, 2(2):140–150, 2003.
- [12] A Savitzky and M J E Golay. Smoothing and differentiation of data by simplified least squares procedures. *Analytical Chemistry*, 36(8):1627–1639, 1964.
- [13] Y H Yang, M J Buckley, S Dudoit, and T P Speed. Comparison of methods for image analysis on cDNA microarray data. *Journal of Computational and Graphical Statistics*, 11:108–136, 2002.
- [14] C.-D. Mayer and C. A. Glasbey. Statistical methods in microarray gene expression data analysis. In D Husmeier, R Dybowski, and S Roberts, editors, *Probabilistic Modeling in Bioinformatics and Medical Informatics*, Advanced Information and Knowledge Processing, pages 211–238. Springer, London, 2005.
- [15] A C Sauve and T P Speed. Normalization, baseline correction and alignment of high-throughput mass spectrometry data. *Proceedings of the Genomic Signal Processing and Statistics Workshop*, 2004.
- [16] O Kohlbacher, K Reinert, C Gröpl, E Lange, N Pfeifer, O Schulz-Trieglaff, and M Sturm. TOPP—the OpenMS proteomics pipeline. *Bioinformatics*, 23(2):e191–e197, 2007.
- [17] E Lange, C Gröpl, O Schulz-Trieglaff, A Leinenbach, C Huber, and K Reinert. A geometric approach for the alignment of liquid chromatography—mass spectrometry data. *Bioinformatics*, 23(13):i273–i281, 2007.
- [18] M Sturm, A Bertsch, C Gröpl, A Hildebrandt, R Hussong, E Lange, N Pfeifer, O Schulz-Trieglaff, A Zerck, K Reinert, and O Kohlbacher. OpenMS - an open-source software framework for mass spectrometry. *BMC Bioinformatics*, 9(1):163, 2008.

- [19] C Bauer, F Kleinjung, C Smith, M Towers, A Tiss, A Chadt, T Dreja, D Beule, H Al-Hasani, K Reinert, J Schuchhardt, and R Cramer. Biomarker discovery and redundancy reduction towards classification using a multi-factorial maldi-tof ms t2dm mouse model dataset. *BMC Bioinformatics*, 12(1):140, 2011.
- [20] M Morháč. An algorithm for determination of peak regions and baseline elimination in spectroscopic data. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 600(2):478–487, 2009.
- [21] C G Ryan, E Clayton, W L Griffin, S H Sie, and D R Cousens. SNIP, a statistics-sensitive background treatment for the quantitative analysis of {PIXE} spectra in geoscience applications. *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms*, 34(3):396–402, 1988.
- [22] C Yang, Z He, and W Yu. Comparison of public peak detection algorithms for maldi mass spectrometry data analysis. *BMC Bioinformatics*, 10(1), 2009.
- [23] E Dougherty. *Mathematical Morphology in Image Processing*. Marcel-Dekker, New York, 1992.
- [24] P Soille. *Morphological Image Analysis: Principles and Applications*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 1999.
- [25] J Y Gil and R Kimmel. Efficient dilation, erosion, opening, and closing algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(12):1606–1617, 2002.
- [26] M Van Droogenbroeck and M J Buckley. Morphological erosions and openings: fast algorithms based on anchors. *Journal of Mathematical Imaging and Vision*, 22(2-3):121–142, 2005.
- [27] M van Herk. A fast algorithm for local minimum and maximum filters on rectangular and octagonal kernels. *Pattern Recognition Letters*, 13(7):517–521, 1992.
- [28] J Gil and M Werman. Computing 2-D min, median, and max filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15:504–507, 1993.

- [29] M van Herk, J C de Munck, J V Lebesque, S Muller, C Rasch, and A Touw. Automatic registration of pelvic computed tomography data and magnetic resonance scans including a full circle method for quantitative accuracy evaluation. *Medical physics*, 25:2054, 1998.
- [30] C Heneghan, J Flynn, M O’Keefe, and M Cahill. Characterization of changes in blood vessel width and tortuosity in retinopathy of prematurity using image analysis. *Medical Image Analysis*, 6(4):407–429, 2002.
- [31] G Zhang, B M Ueberheide, S Waldemarson, S Myung, K Molloy, J Eriksson, B T Chait, T A Neubert, and D Fenyö. Protein quantitation using mass spectrometry. In D Fenyö, editor, *Computational Biology*, volume 673 of *Methods in Molecular Biology*, pages 211–222. Humana Press, New York, 2010.
- [32] L Greengard and J-Y Lee. Accelerating the nonuniform fast fourier transform. *SIAM review*, 46(3):443–454, 2004.
- [33] A W Lo and A C MacKinlay. An econometric analysis of nonsynchronous trading. *Journal of Econometrics*, 45(1–2):181–211, 1990.
- [34] A Aris, B Shneiderman, C Plaisant, G Shmueli, and W Jank. Representing unevenly-spaced time series data for visualization and interactive exploration. In *Human-Computer Interaction-INTERACT 2005*, pages 835–846. Springer, Heidelberg, 2005.
- [35] M Schulz and M Mudelsee. REDFIT: estimating red-noise spectra directly from unevenly spaced paleoclimatic time series. *Computers & Geosciences*, 28(3):421–426, 2002.
- [36] T J Deeming. Fourier analysis with unequally-spaced data. *Astrophysics and Space Science*, 36(1):137–158, 1975.
- [37] J D Scargle. Studies in astronomical time series analysis. II-Statistical aspects of spectral analysis of unevenly spaced data. *The Astrophysical Journal*, 263:835–853, 1982.
- [38] M Bourgeois, FTAW Wajer, D van Ormondt, and D Graveron-Demilly. Reconstruction of MRI images from non-uniform sampling and its application to intrascan motion correction in functional MRI. In J J Benedetto and P J S G Ferreira, editors, *Modern Sampling Theory*, Applied and Numerical Harmonic Analysis, pages 343–363. Birkhäuser, Boston, 2001.

- [39] G M Fiedler, S Baumann, A Leichtle, A Oltmann, J Kase, J Thiery, and U Ceglarek. Standardized peptidome profiling of human urine by magnetic bead separation and matrix-assisted laser desorption/ionization time-of-flight mass spectrometry. *Clinical Chemistry*, 53(3):421–428, 2007.
- [40] P B Yildiz, Y Shyr, J S M Rahman, N R Wardwell, L J Zimmerman, B Shakh-tour, W H Gray, S Chen, M Li, H Roder, D C Liebler, W L Bigbee, J M Siegfried, J L Weissfeld, A L Gonzalez, M Ninan, D H Johnson, D P Carbone, R M Caprioli, and P P Massion. Diagnostic accuracy of MALDI mass spectrometric analysis of unfractionated serum in lung cancer. *Journal of thoracic oncology: official publication of the International Association for the Study of Lung Cancer*, 2(10):893, 2007.
- [41] B Wu, T Abbott, D Fishman, W McMurray, G Mor, K Stone, D Ward, K Williams, and H Zhao. Comparison of statistical methods for classification of ovarian cancer using mass spectrometry data. *Bioinformatics*, 19(13):1636–1643, 2003.
- [42] W Yu, X Li, J Liu, B Wu, K R Williams, and H Zhao. Multiple peak alignment in sequential data analysis: a scale-space-based approach. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 3(3):208–219, 2006.
- [43] B-L Adam, Y Qu, J W Davis, M D Ward, O J Semmes, P F Schellhammer, Y Yasui, Z Feng, G L Wright Jr., M A Clements, and L H Cazares. Serum protein fingerprinting coupled with a pattern-matching algorithm distinguishes prostate cancer from benign prostate hyperplasia and healthy men. *Cancer Research*, 62(13):3609–3614, 2002.
- [44] F Taguchi, B Solomon, V Gregorc, H Roder, R Gray, K Kasahara, M Nishio, J Brahmer, A Spreafico, V Ludovini, P P Massion, R Dziadziuszko, J Schiller, J Grigorieva, M Tsypin, S W Hunsucker, R Caprioli, M W Duncan, F R Hirsch, P A Bunn, and D P Carbone. Mass spectrometry to classify non-small-cell lung cancer patients for clinical outcome after treatment with epidermal growth factor receptor tyrosine kinase inhibitors: A multicohort cross-institutional study. *Journal of the National Cancer Institute*, 99(11):838–846, 2007.
- [45] M Li, S Chen, J Zhang, H Chen, and Y Shyr. Wave-spec: a preprocessing package for mass spectrometry data. *Bioinformatics*, 27(5):739–740, 2011.

- [46] D Mantini, F Petrucci, D Pieragostino, P Del Boccio, M Di Nicola, C Di Ilio, G Federici, P Sacchetta, S Comani, and A Urbani. LIMPIC: a computational method for the separation of protein MALDI-TOF-MS signals from noise. *BMC Bioinformatics*, 8(1):101, 2007.
- [47] G Siuzdak. *The Expanding Role of Mass Spectrometry in Biotechnology*. McC Press, 2006.
- [48] L L House, M A Clyde, and R L Wolpert. Bayesian nonparametric models for peak identification in MALDI-TOF mass spectroscopy. *The Annals of Applied Statistics*, 5(2B):1488–1511, 2011.
- [49] R J Hyndman and A B Koehler. Another look at measures of forecast accuracy. *International Journal of Forecasting*, 22(4):679–688, 2006.
- [50] A Eckner. Algorithms for unevenly-spaced time series: Moving averages and other rolling operators. Technical report, Working Paper, 2013. http://www.eckner.com/papers/ts_alg.pdf.

Appendix A Six transformation functions on each of the six datasets

The effect of the six transformation functions, t_0 - t_5 , on a spectrum from each of the six datasets are shown in Figures A.1 to A.6.

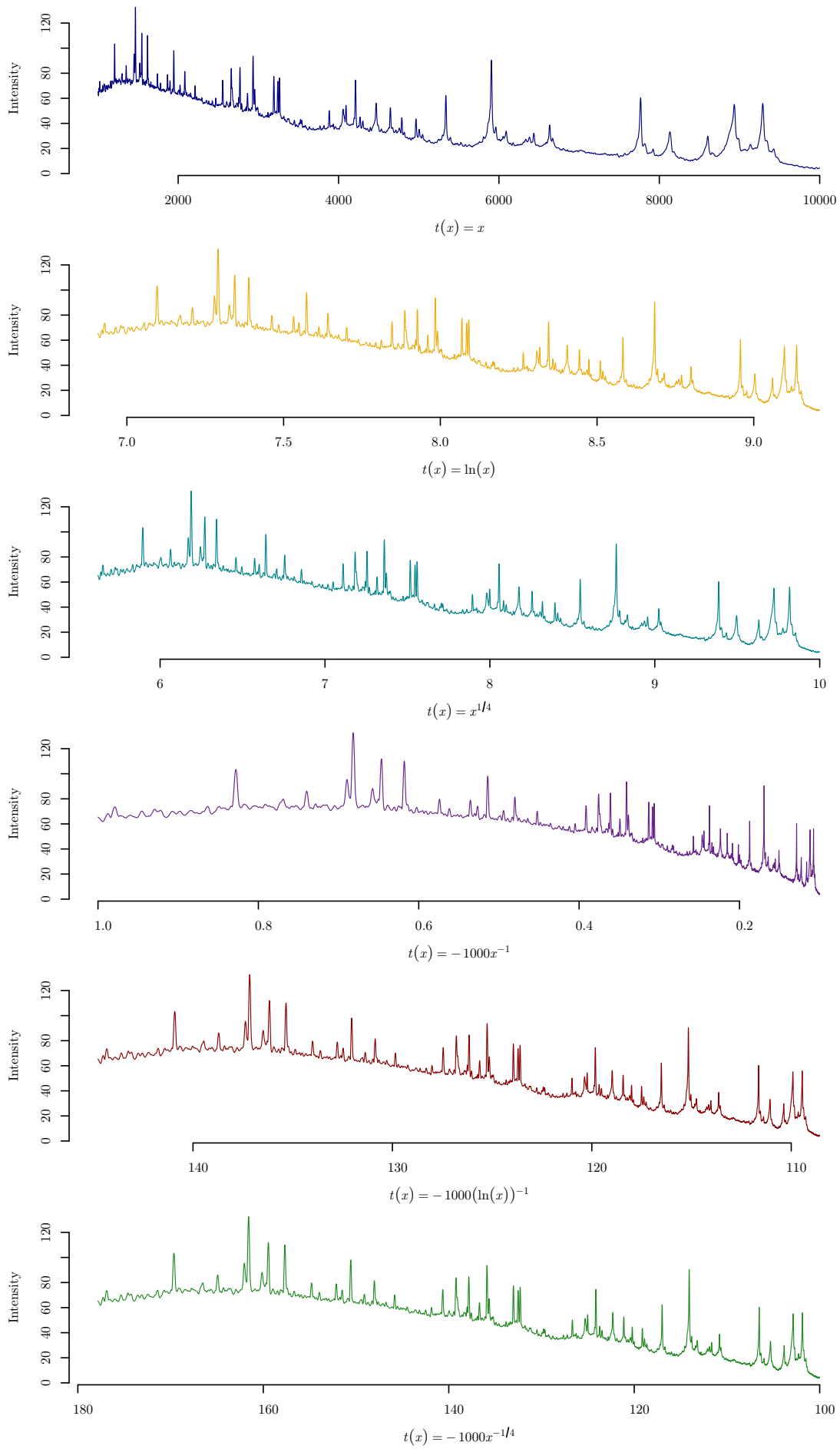


Figure A.1: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Fiedler dataset.

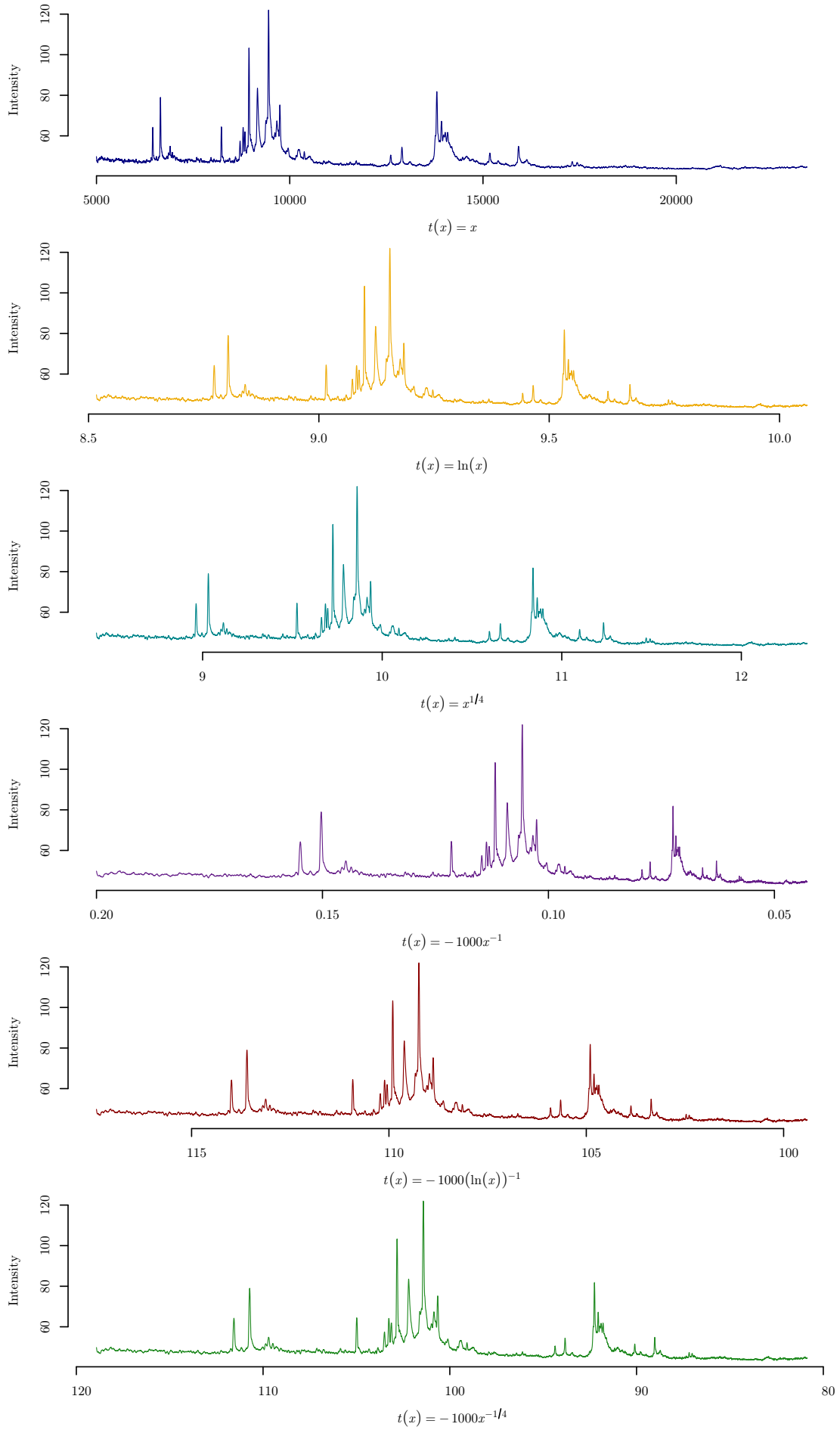


Figure A.2: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Yildiz dataset.

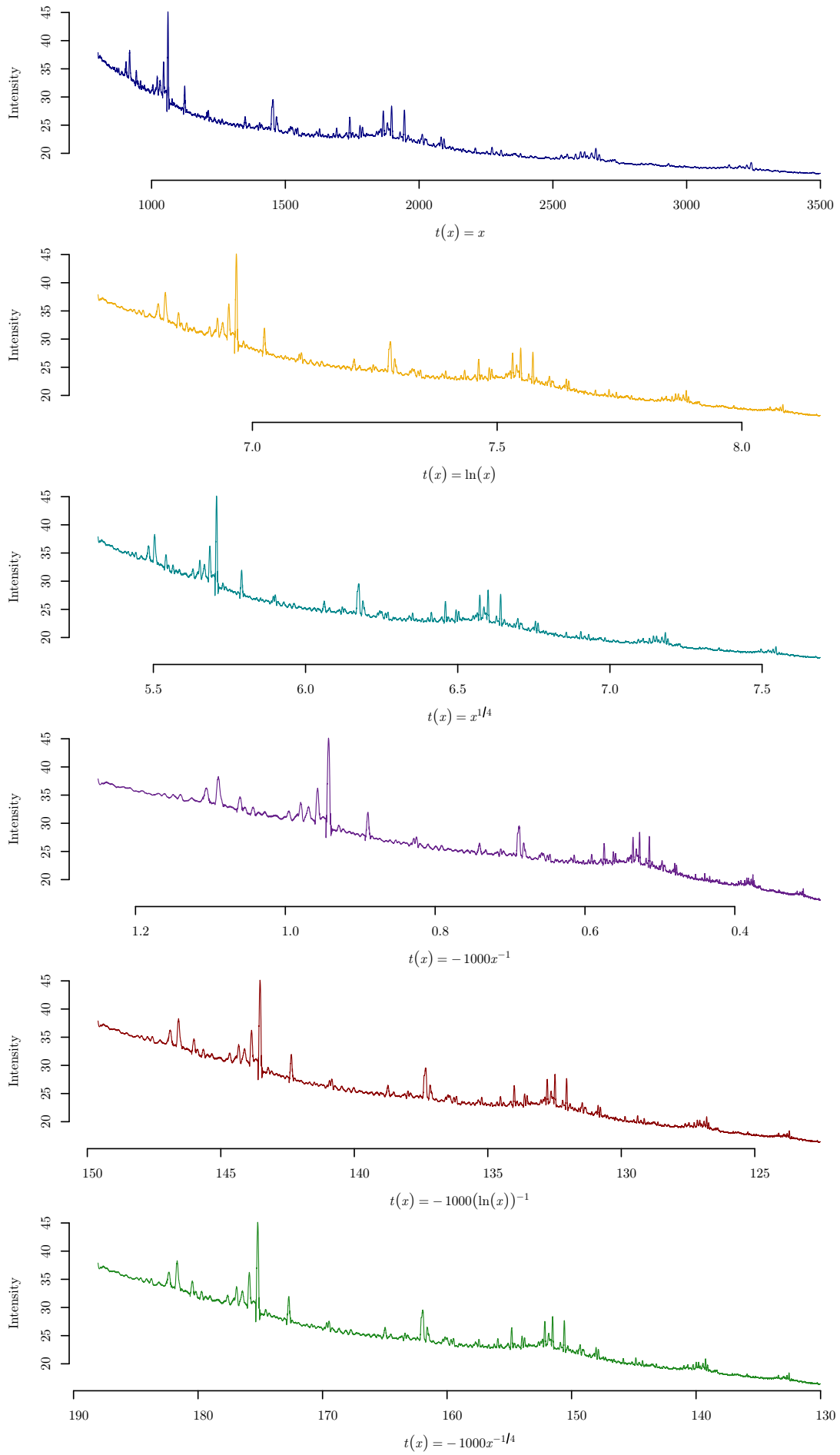


Figure A.3: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Wu dataset.

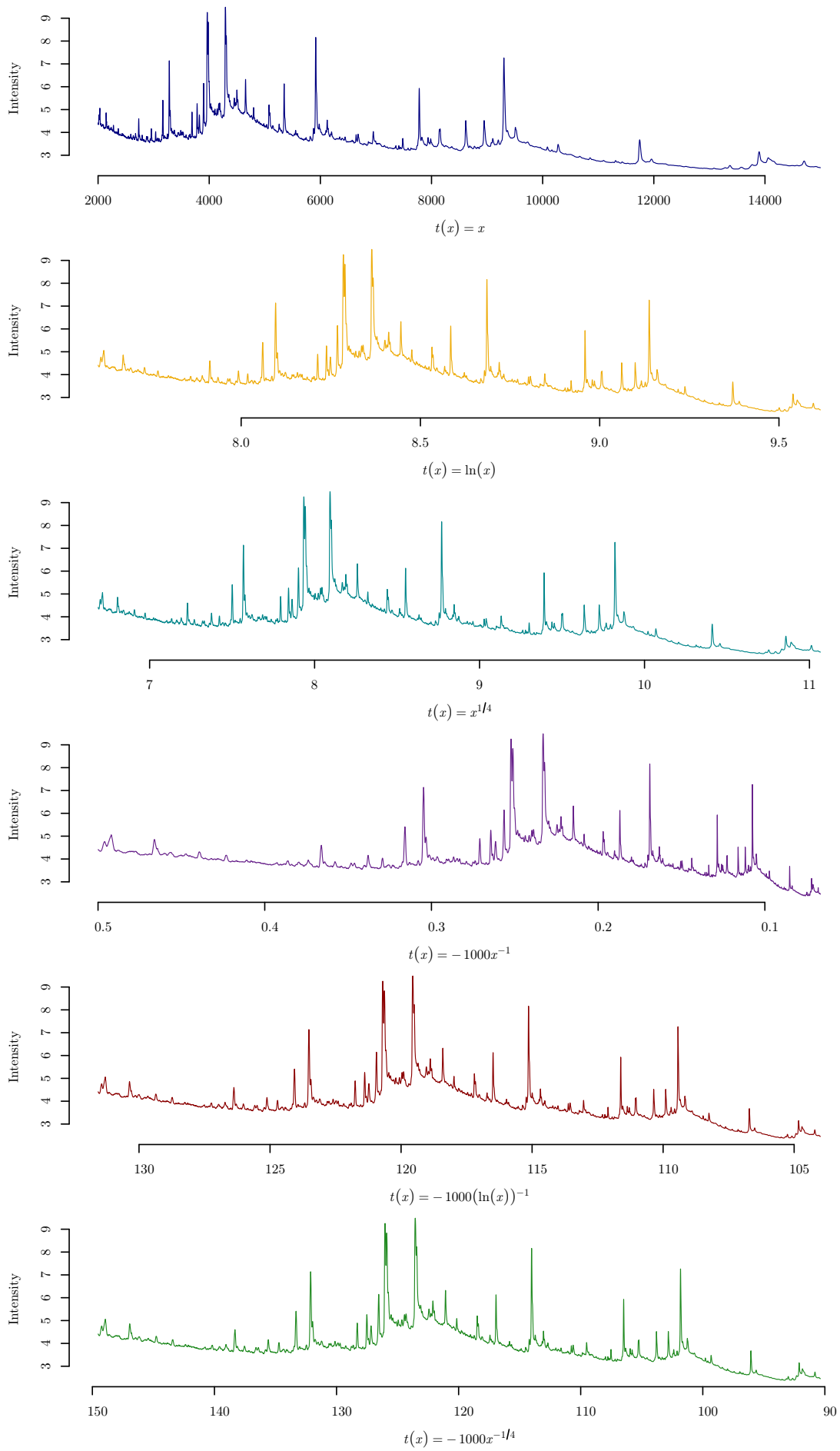


Figure A.4: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Adam dataset.

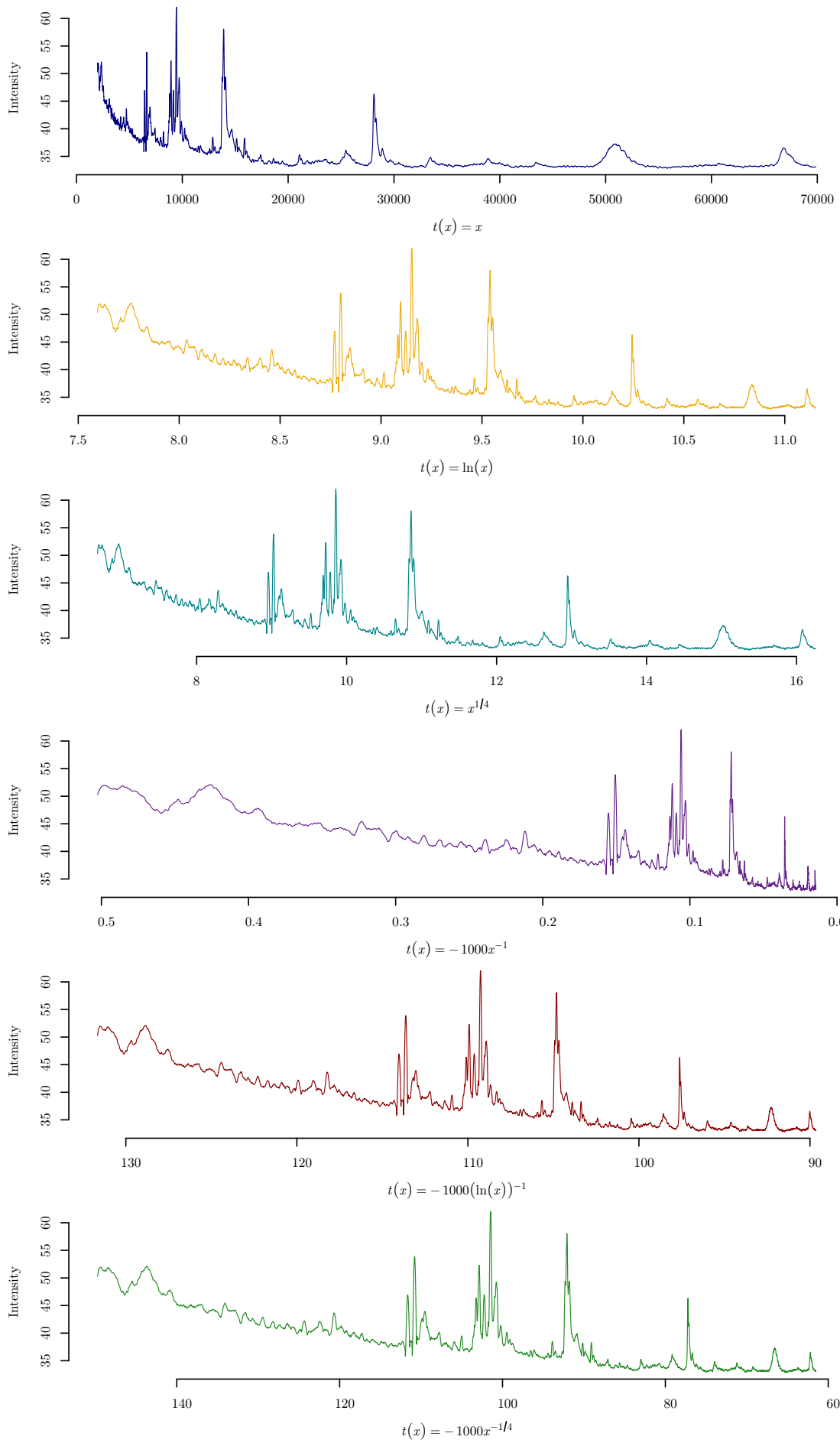


Figure A.5: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Taguchi dataset.

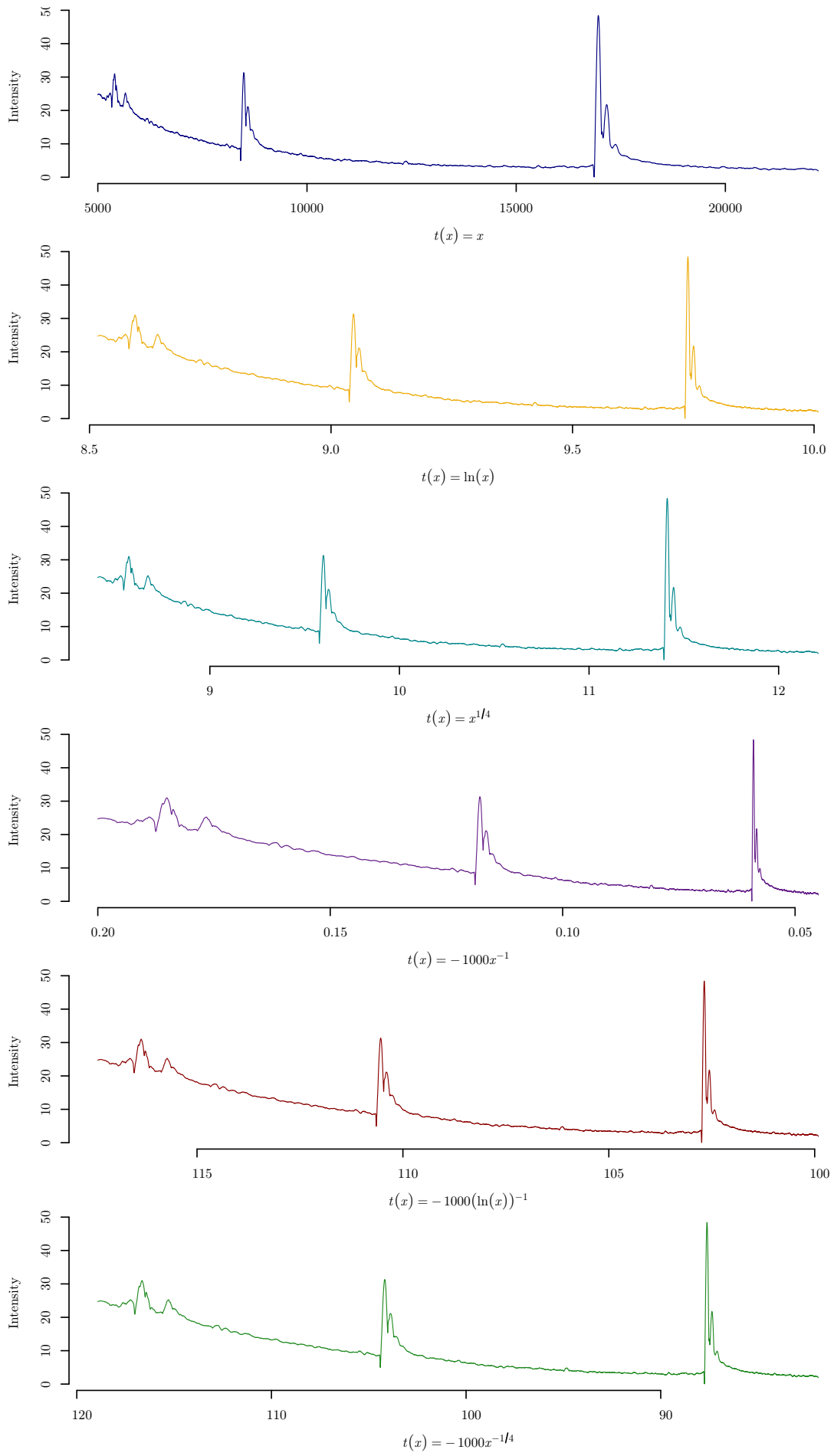


Figure A.6: The effect of the six transformation functions, t_0 - t_5 , on a spectrum from the Mantini dataset.

Appendix B CLSA examples

To illustrate how the CLSA works, consider two cases of the algorithm in returning the erosion in Figures B.1 and B.2.

Figure B.1 shows a case where $\epsilon_B(f)(x_{12}) = 3$ using a SE of size $k = 3$. It can be seen that,

$$\theta_{w_{12}^{\nabla}-1} = \theta_{10} = 3 \neq \theta_{w_{12}^{\Delta}} = \theta_{13} = 4,$$

but,

$$\theta_{w_{12}^{\nabla}} = \theta_{11} = 3 = \theta_{w_{12}^{\Delta}+1} = \theta_{14}.$$

Therefore the desired result is also achieved using the CLSA as,

$$r_{\min}(f(x_{12})) = g(x_{w_{12}^{\Delta}}) = g(x_{13}) = 3 = \epsilon_B(f)(x_{12}).$$

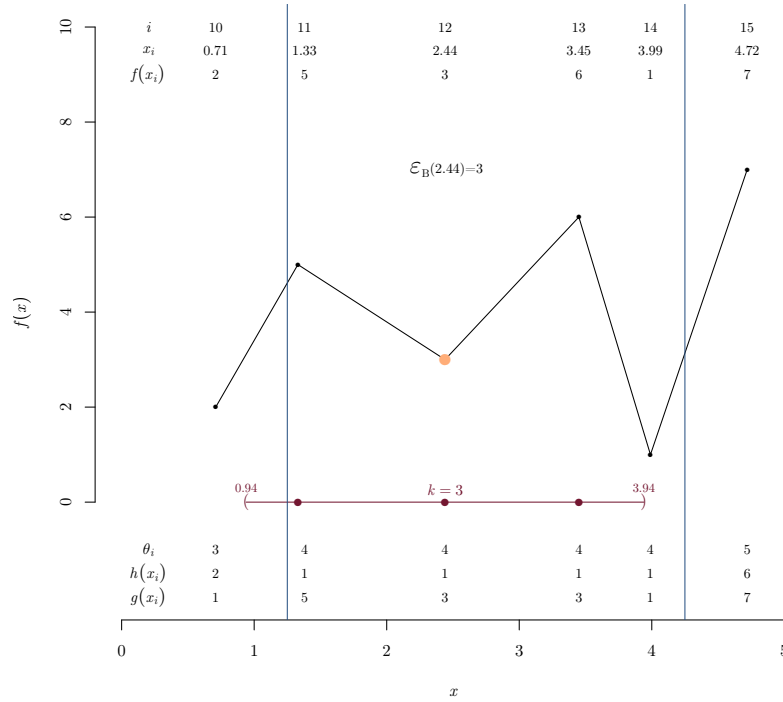


Figure B.1: An example of data for $x_i = x_{12} = 2.44$ and $k = 3$ where $\theta_{w_i^{\nabla}} = \theta_{w_i^{\Delta}+1}$ (i.e. $\theta_{w_{12}^{\nabla}} = \theta_{w_{12}^{\Delta}+1} = 4$) and the computation required to return the result of the CLSA (the tan coloured point $f(2.44) = 3$).

Figure B.2 is a different case in the CLSA where $\theta_{w_i^\nabla-1} = \theta_{w_i^\Delta}$, as opposed to the case shown in Figure B.1 where $\theta_{w_i^\nabla} = \theta_{w_i^\Delta+1}$. To obtain the erosion of point $x_i = x_9 = 2.44$ for $k = 3$ using the CLSA, observe that

$$\theta_{w_9^\nabla} = \theta_8 = 3 \neq \theta_{w_9^\Delta+1} = \theta_{11} = 4,$$

and

$$\theta_{w_9^\nabla-1} = \theta_7 = 3 = \theta_{w_9^\Delta} = \theta_{10}.$$

Therefore, the result of the CLSA erosion is

$$r_{\min}(f(x_9)) = h(x_{w_9^\nabla}) = h(x_8) = 3.$$

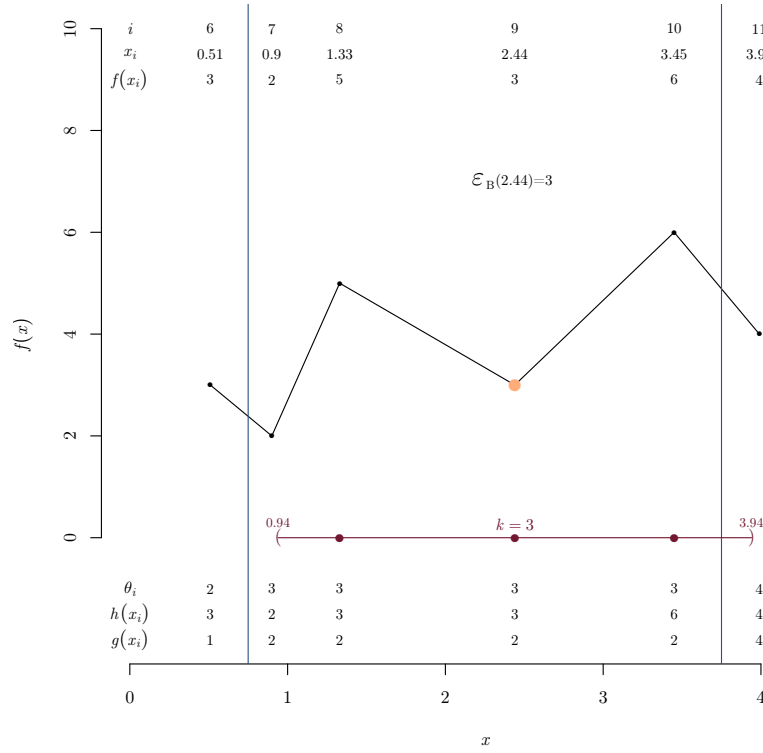


Figure B.2: An example of data where $\theta_{w_i^\nabla-1} = \theta_{w_i^\Delta}$ and the computation required for the continuous line segment algorithm.

Appendix C CLSA and naive algorithm computational times for simulated data

The R code required to produce the top-hat baseline subtraction computational time results shown in Table 4 is presented below.

```
#####
#####
##### preamble and set-up #####
#####
#####

# required to install github packages
# see http://cran.r-project.org/web/packages/devtools/README.html
library(devtools)

#install CLSA package
devtools::install_github('tystan/clsa')
library(clsa)
# documentation
?clsa_min

# latex formatted table output
library(xtable)

# function to create x values (m/z values)
get_x_coords<-function(n) return(sort(rbeta(n,1,3)))
# function to create intensity values
get_f_signal<-function(n) return(rchisq(n,10))

#####
#####
##### testing comp times #####
#####
#####

# creating a data.frame with the following columns:
# * n: the number of m/z points
# * win: the window sizes
# * time_naiv: using the naive alg -- time taken for the row's 'n' and 'win' values
# * time_clsa: using the CLSA -- time taken for the row's 'n' and 'win' values

# ranges of dataset size and window size:
```



```

n_rng<-seq(1e4,1e5,by=1e4)
win_rng<-c(0.005,0.01,0.02,0.05,0.1,0.2)
# these times will be updated
time_naiv<-0
time_clsa<-0
# enumerate all n and win combinations
time_df<-as.data.frame(
  expand.grid(
    n=n_rng
    ,win=win_rng
    ,time_naiv=time_naiv
    ,time_clsa=time_clsa
  )
)
(n_df<-nrow(time_df)) # should be 6x10=60
time_df

# test each n and win combination 20 times,
# i.e., each dataset has 20 "spectra"
n_reps<-20
set.seed(12345) # make reproducible

# now iterate over rows of the data.frame for each "spectrum",
# and time the computation
for(j in 1:n_reps)
{
  for(i in 1:n_df)
  {
    # get n and win combination
    n<-time_df$n[i]
    this_win<-time_df$win[i]
    # update console of progress
    cat("::: Iteration", (j-1)*n_df+i, "of", n_df*n_reps, "::: ")
    cat("n =", n, "and window size = ", this_win, ":\n")
    # randomly generate x and f
    x<-get_x_coords(n)
    f<-get_f_signal(n)

    # time the computations, add to previous "spectra" times
    time_df$time_clsa[i]<-time_df$time_clsa[i]+
      system.time(a<-clsa_max(x,clsa_min(x,f,this_win),this_win))[3]
    time_df$time_naiv[i]<-time_df$time_naiv[i]+
      system.time(b<-naiv_max(x,naiv_min(x,f,this_win),this_win))[3]

    # if the results are not equal between the naiv and CLSA we have

```

```

    # a problem; ABORT!
    if(!all(a==b)) break;
  }
}
time_df

#####
##### table of results #####
#####

# extract times for naiv and CLSA for printing
time_naiv<-data.frame(n=time_df$n,win=time_df$win,time=time_df$time_naiv,func="Naive")
time_clsa<-data.frame(n=time_df$n,win=time_df$win,time=time_df$time_clsa,func="CLSA")

# print the results!
xtable(
  cbind(
    xtabs(time ~ I(n/1e4) + win, data = time_naiv)
    ,NA
    ,xtabs(time ~ I(n/1e4) + win, data = time_clsa)
  )
  ,digits=1
)

```